

# Data thinning to avoid double dipping

Anna Neufeld  
BIRS Workshop  
February, 2024

## What is double dipping?

---

Classical statistical methods assume that we only ever test pre-specified hypotheses about pre-specified models.

# What is double dipping?

---

Classical statistical methods assume that we only ever test pre-specified hypotheses about pre-specified models.

In reality, we explore our data, fit several models, evaluate these models, select our favorite model, then test hypotheses about this model.

# What is double dipping?

---

Classical statistical methods assume that we only ever test pre-specified hypotheses about pre-specified models.

In reality, we explore our data, fit several models, evaluate these models, select our favorite model, then test hypotheses about this model.

**Double Dipping:** Using the same data for two tasks, such as:

1. Generating and testing a null hypothesis.
2. Fitting and evaluating a model.

# One approach: avoid double dipping entirely via sample splitting

---

	<b>Feature 1</b>	<b>Feature 2</b>
<b>Obs. 1</b>	12	6
<b>Obs. 2</b>	31	8
<b>Obs. 3</b>	11	31
<b>Obs. 4</b>	22	34

# One approach: avoid double dipping entirely via sample splitting

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8
Obs. 3	11	31
Obs. 4	22	34

Train

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8

Test

	Feature 1	Feature 2
Obs. 3	11	31
Obs. 4	22	34

# One approach: avoid double dipping entirely via sample splitting

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8
Obs. 3	11	31
Obs. 4	22	34

Train

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8

Select hypothesis.

Test

	Feature 1	Feature 2
Obs. 3	11	31
Obs. 4	22	34

# One approach: avoid double dipping entirely via sample splitting

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8
Obs. 3	11	31
Obs. 4	22	34

Train

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8

Select hypothesis.

Test

	Feature 1	Feature 2
Obs. 3	11	31
Obs. 4	22	34

Test hypothesis.

# One approach: avoid double dipping entirely via sample splitting

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8
Obs. 3	11	31
Obs. 4	22	34

Train

	Feature 1	Feature 2
Obs. 1	12	6
Obs. 2	31	8

Test

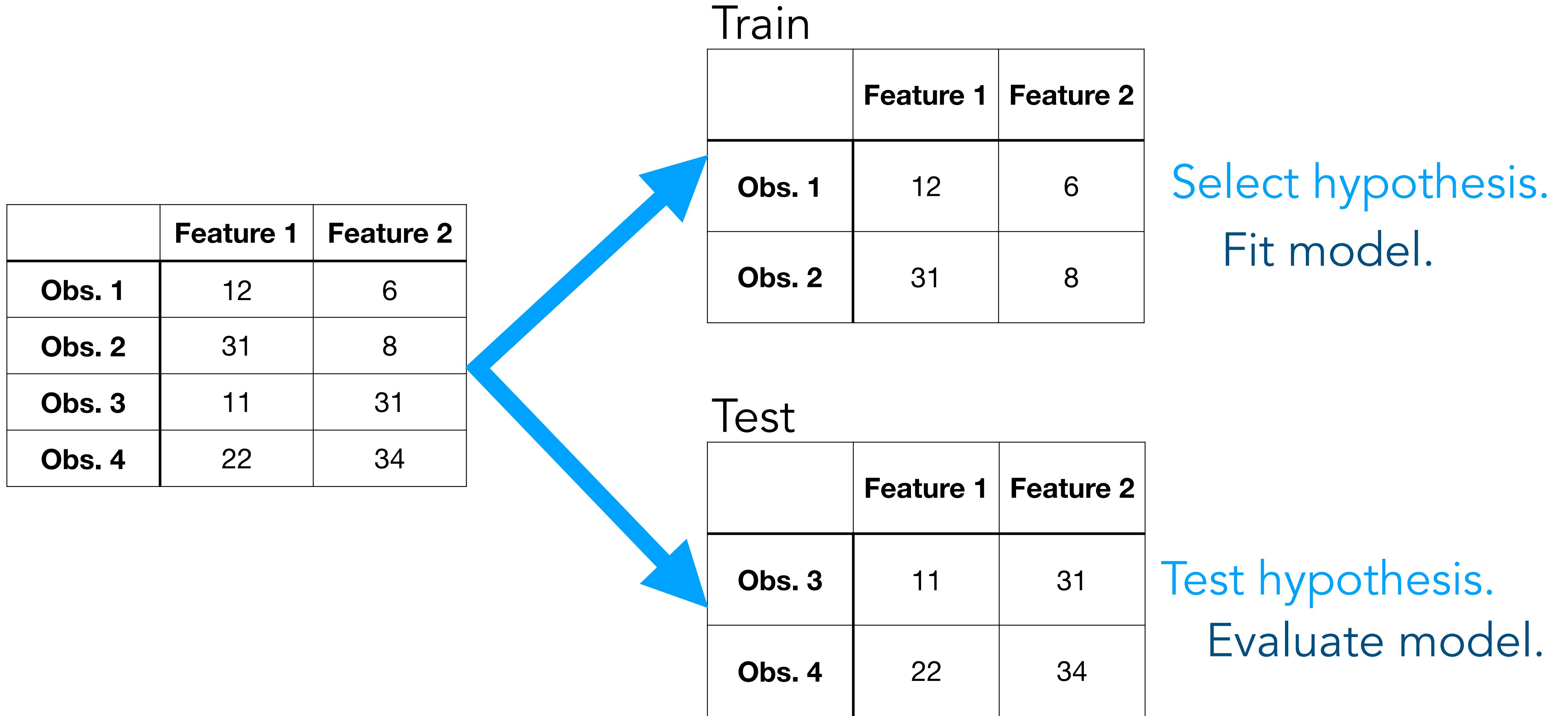
	Feature 1	Feature 2
Obs. 3	11	31
Obs. 4	22	34

Select hypothesis.

Fit model.

Test hypothesis.

# One approach: avoid double dipping entirely via sample splitting



# Outline

---

- 1. Motivation: settings where sample splitting doesn't work**
2. Poisson thinning
3. Data thinning
4. Application to single-cell RNA sequencing data
5. Ongoing work

# Motivating example: identifying differentially expressed genes in scRNA-seq data

---

## scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

# Motivating example: identifying differentially expressed genes in scRNA-seq data

---

## scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

**Question:** Which genes are differentially expressed across cell type?

# Motivating example: identifying differentially expressed genes in scRNA-seq data

---

## scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

**Question:** Which genes are differentially expressed across cell type?

**Problem:** Cell type is unobserved, and must be estimated via clustering.

# Motivating example: identifying differentially expressed genes in scRNA-seq data

## scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

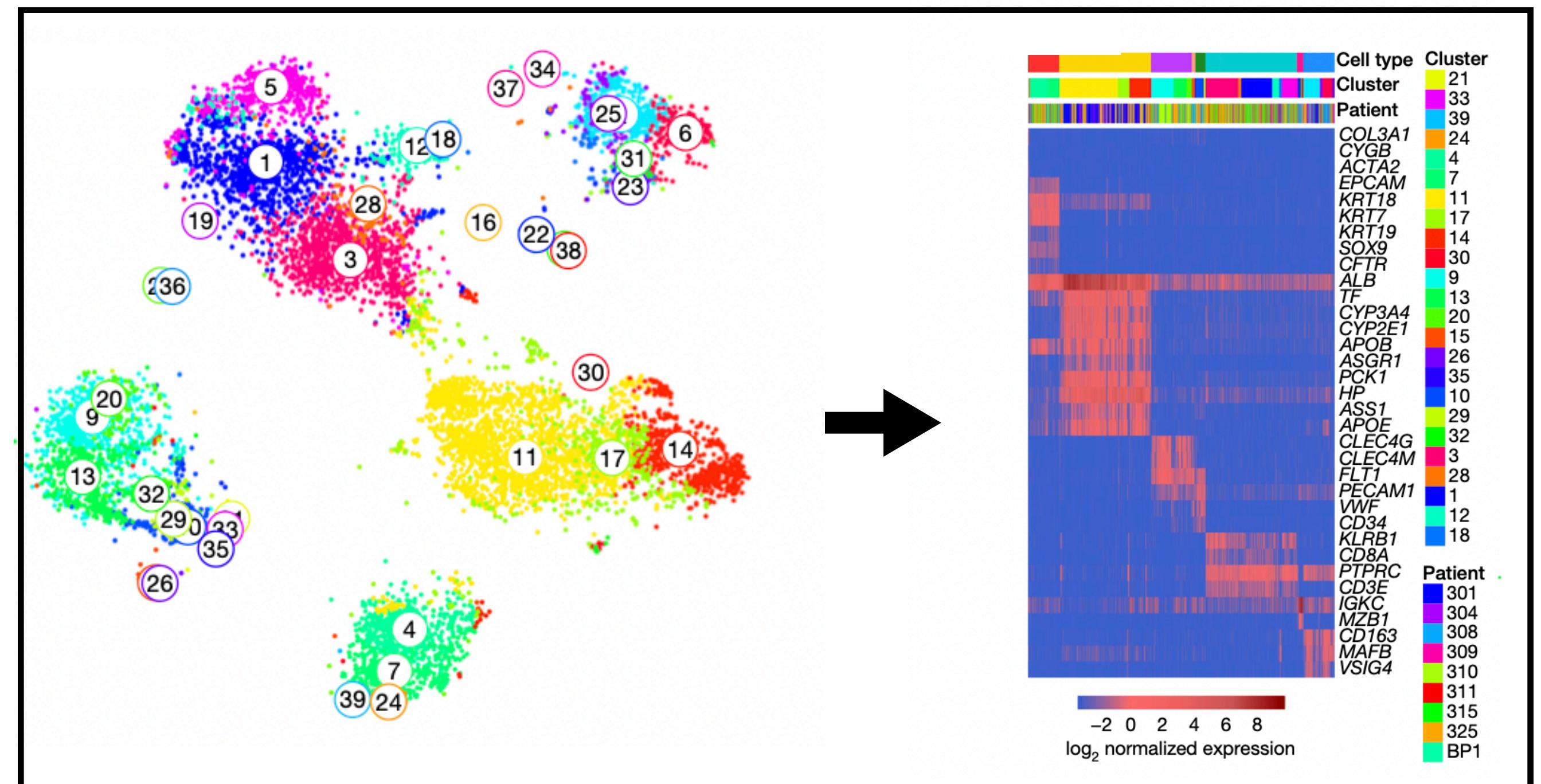
**Question:** Which genes are differentially expressed across cell type?

**Problem:** Cell type is unobserved, and must be estimated via clustering.

**A human liver cell atlas reveals heterogeneity and epithelial progenitors**

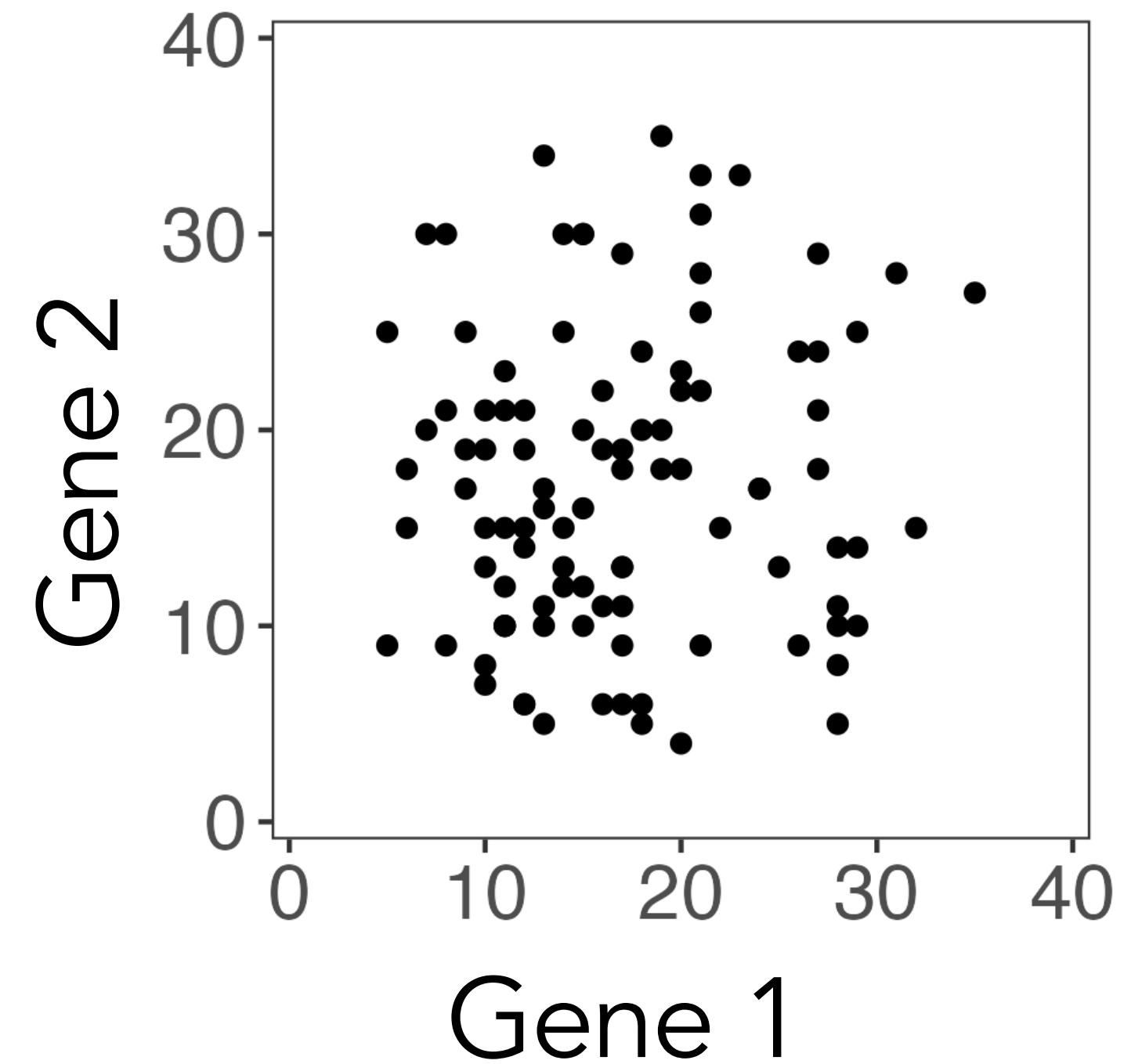
Nadim Aizarani, Antonio Saviano, Sagar, Laurent Mailly, Sarah Durand, Josip S. Herman, Patrick Pessaix, Thomas F. Baumert & Dominic Grün

*Nature* 572, 199–204 (2019) | [Cite this article](#)  
64k Accesses | 284 Citations | 321 Altmetric | [Metrics](#)



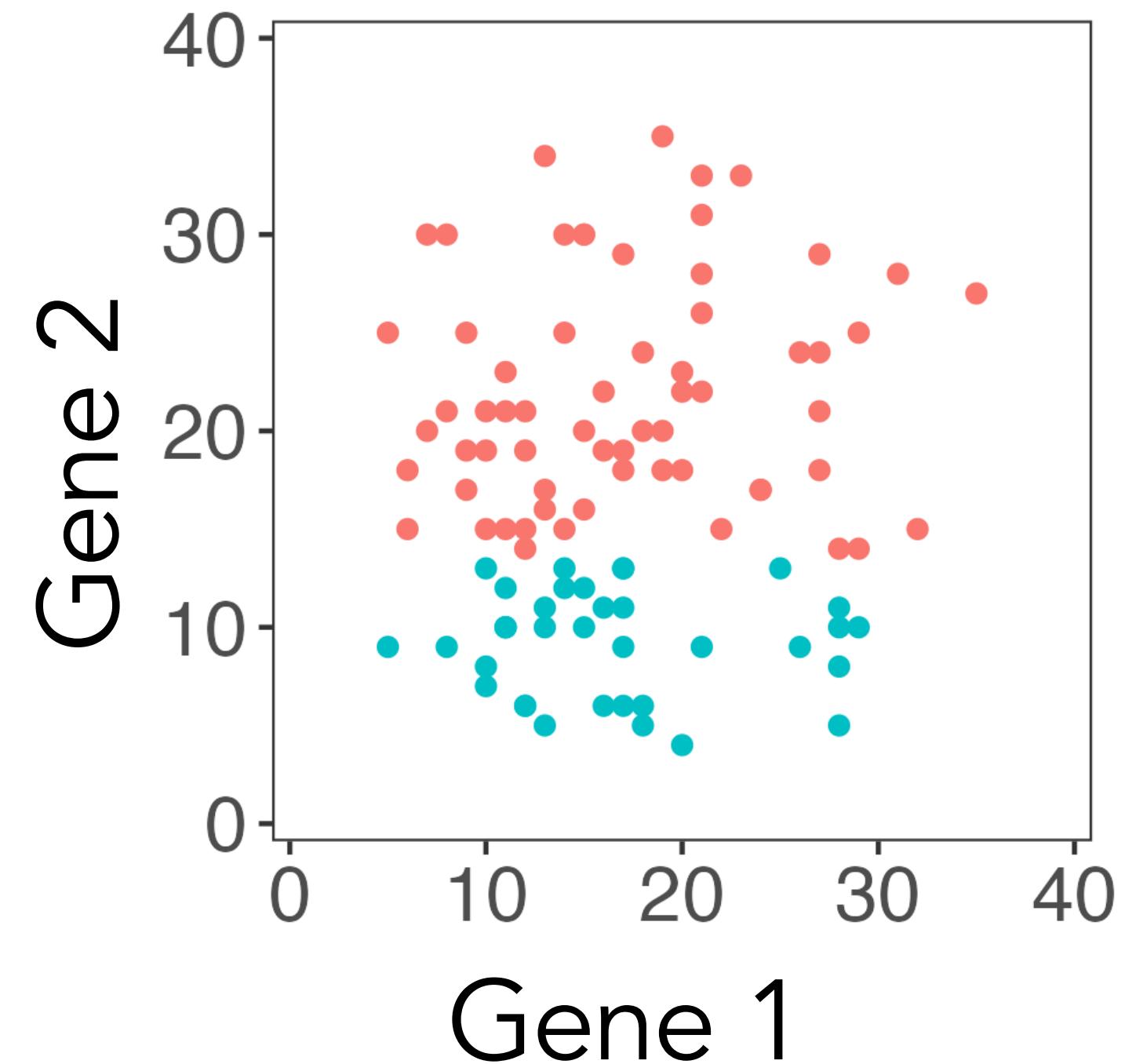
Naively testing for a difference in means across estimated clusters is an example of double dipping

---



Naively testing for a difference in means across estimated clusters is an example of double dipping

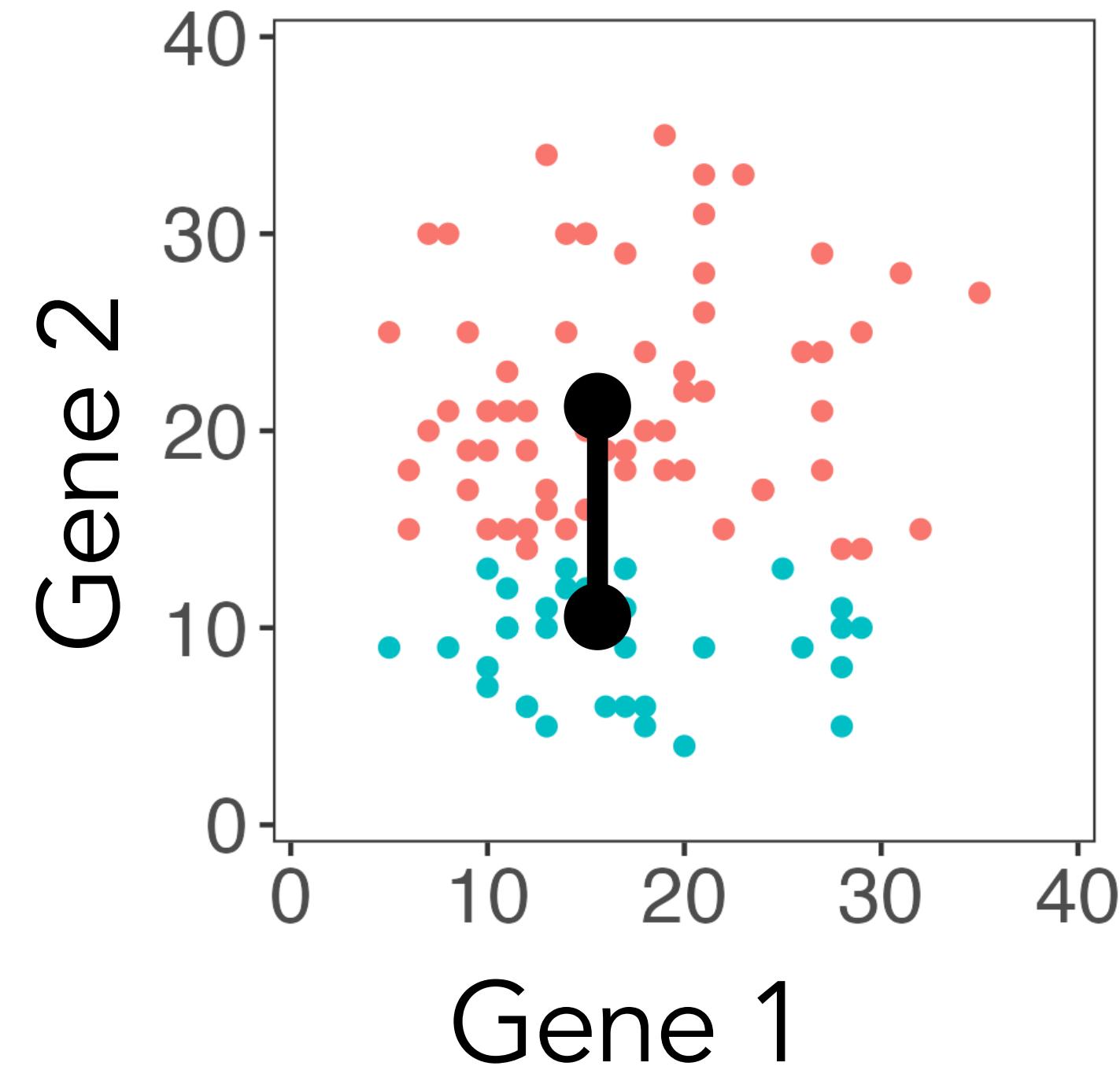
---



**Step 1:** cluster the observations.

Naively testing for a difference in means across estimated clusters is an example of double dipping

---

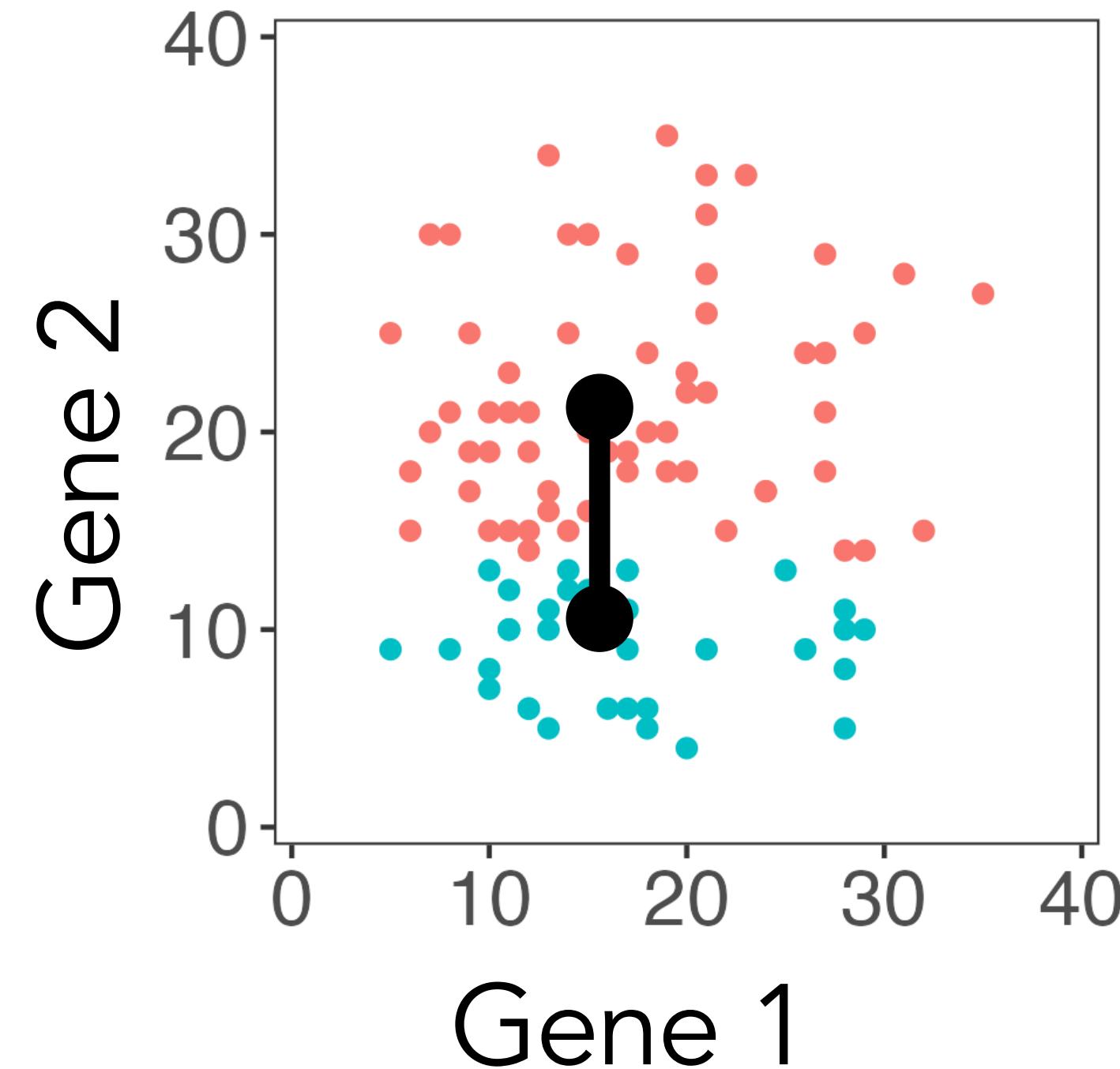


**Step 1:** cluster the observations.

Generate  $H_0$  : “the expected value of Gene 2 is the same between red cells and the blue cells.”

Naively testing for a difference in means across estimated clusters is an example of double dipping

---



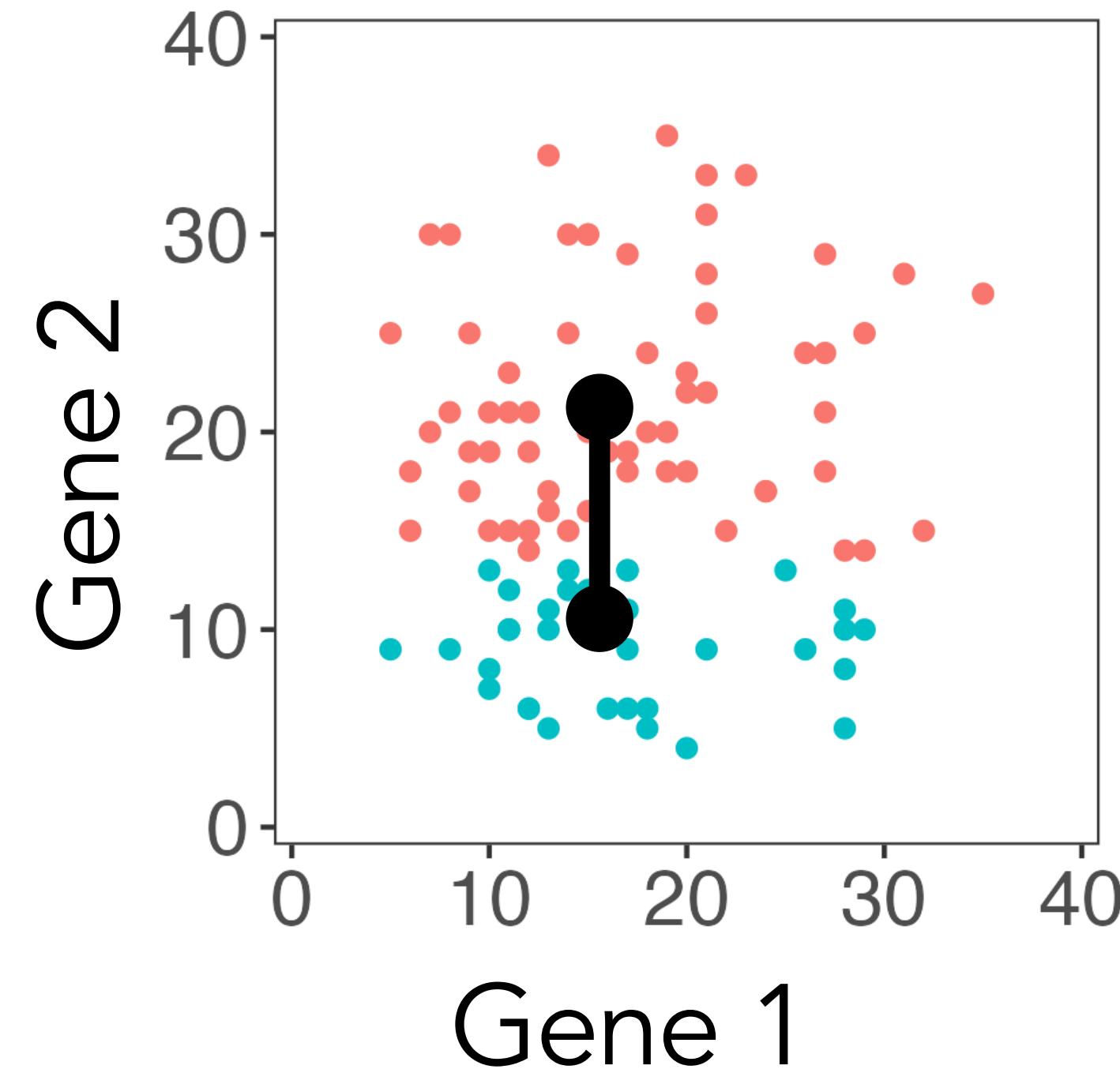
**Step 1:** cluster the observations.

Generate  $H_0$  : “the expected value of Gene 2 is the same between red cells and the blue cells.”

**Step 2:** test  $H_0$  with a t-test.

Naively testing for a difference in means across estimated clusters is an example of double dipping

---



**Step 1:** cluster the observations.

Generate  $H_0$  : “the expected value of Gene 2 is the same between red cells and the blue cells.”

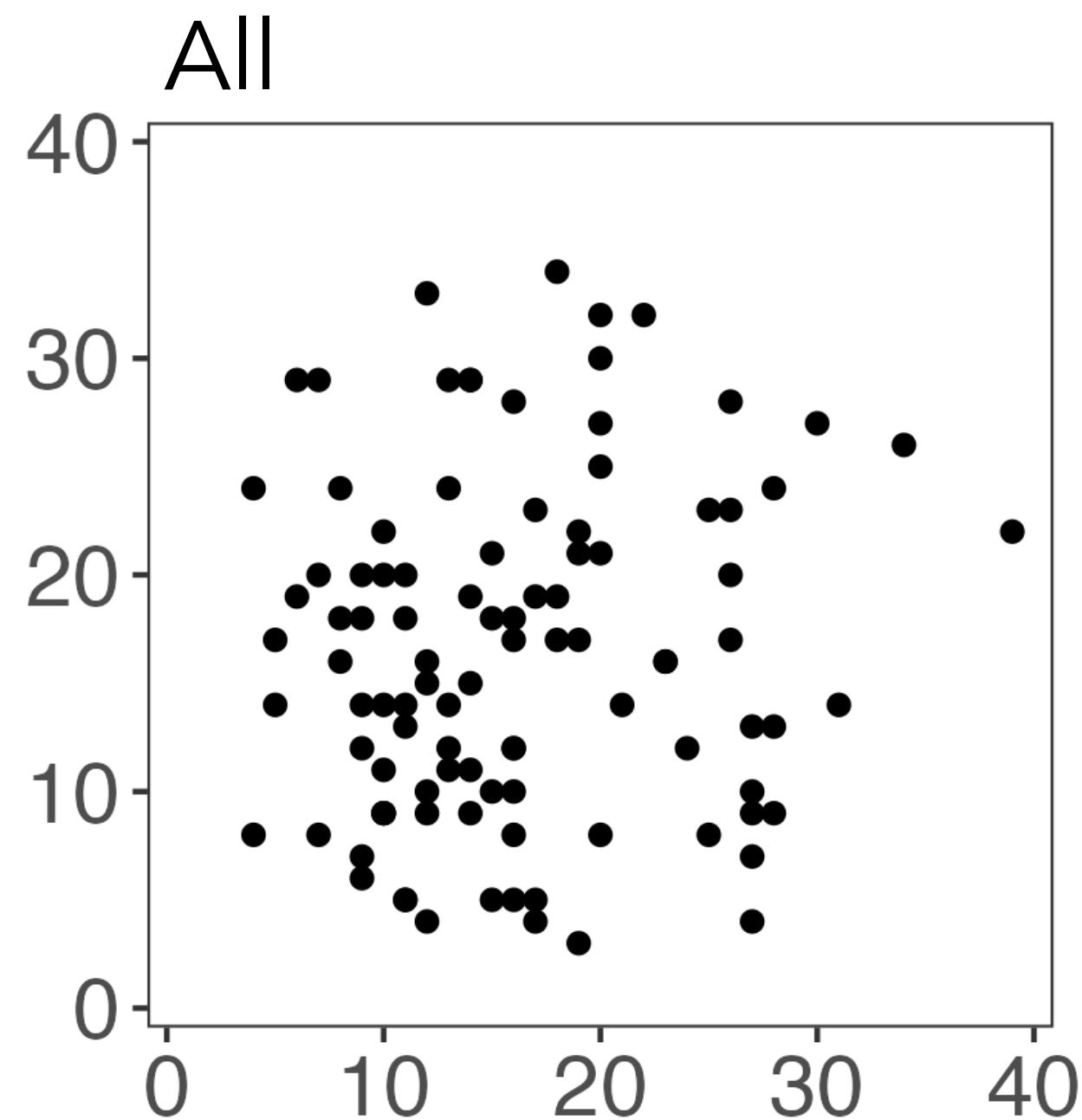
$$p < 10^{-10}$$



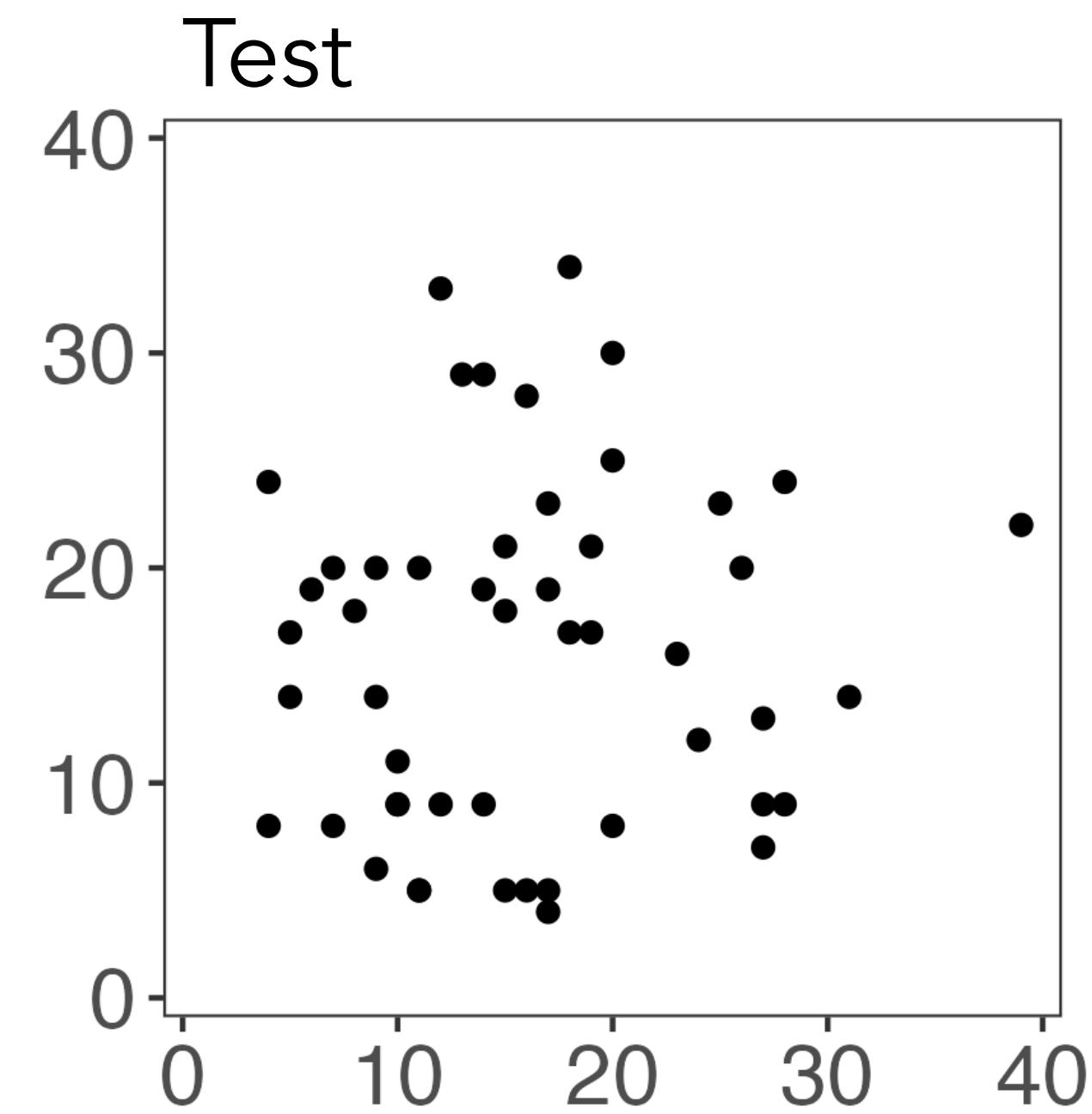
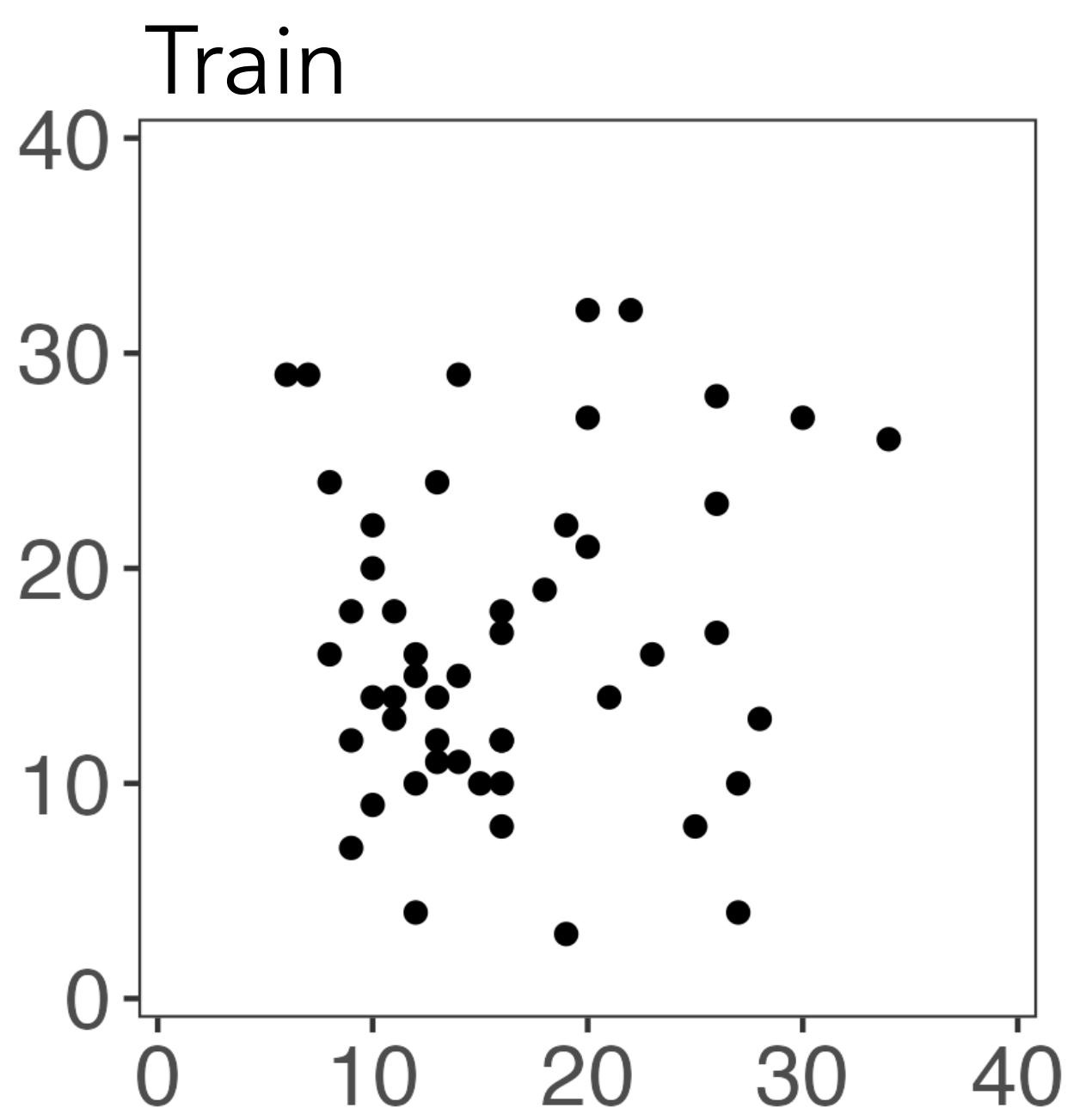
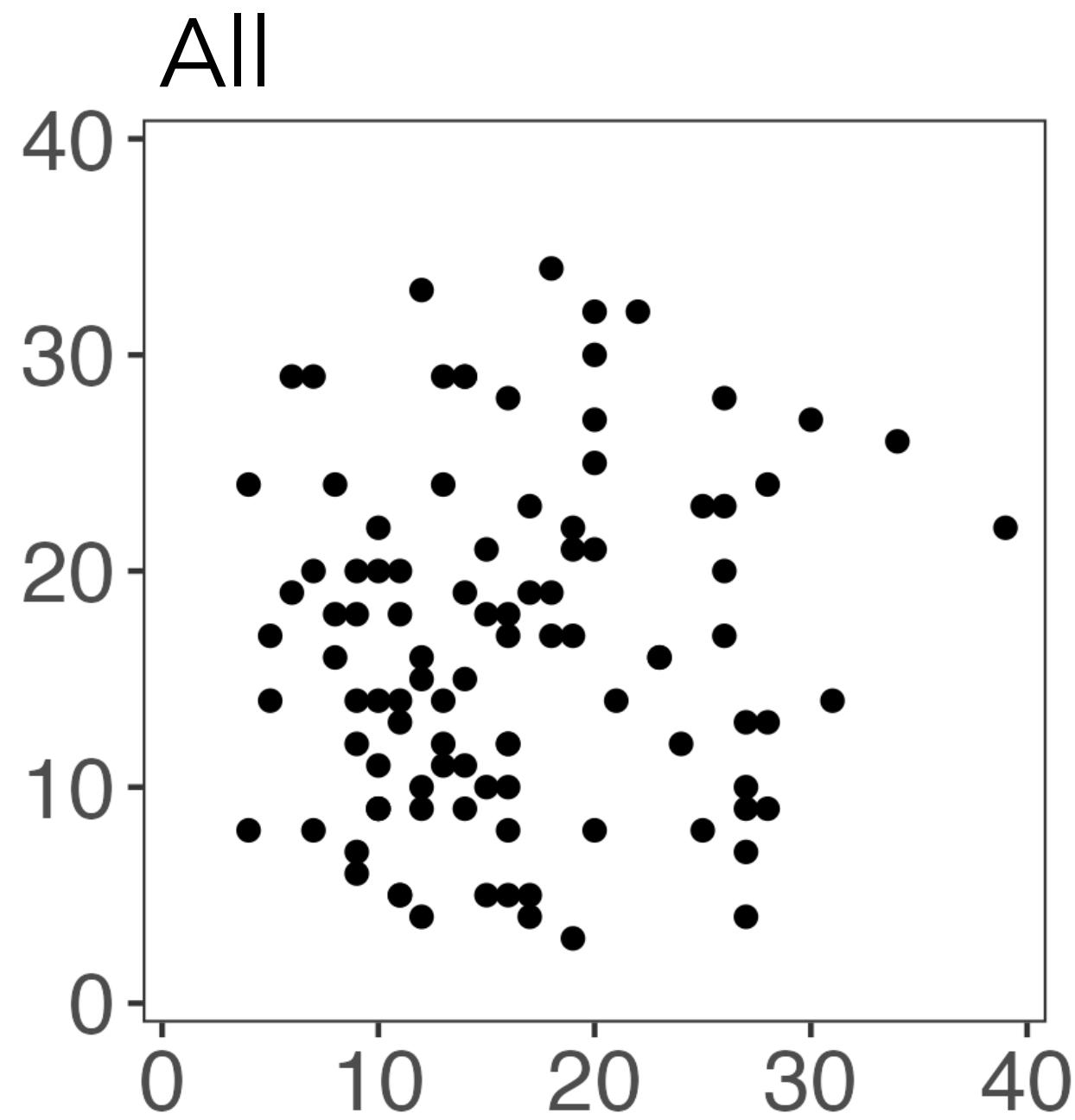
**Step 2:** test  $H_0$  with a t-test.

Sample splitting cannot be used in this context

---

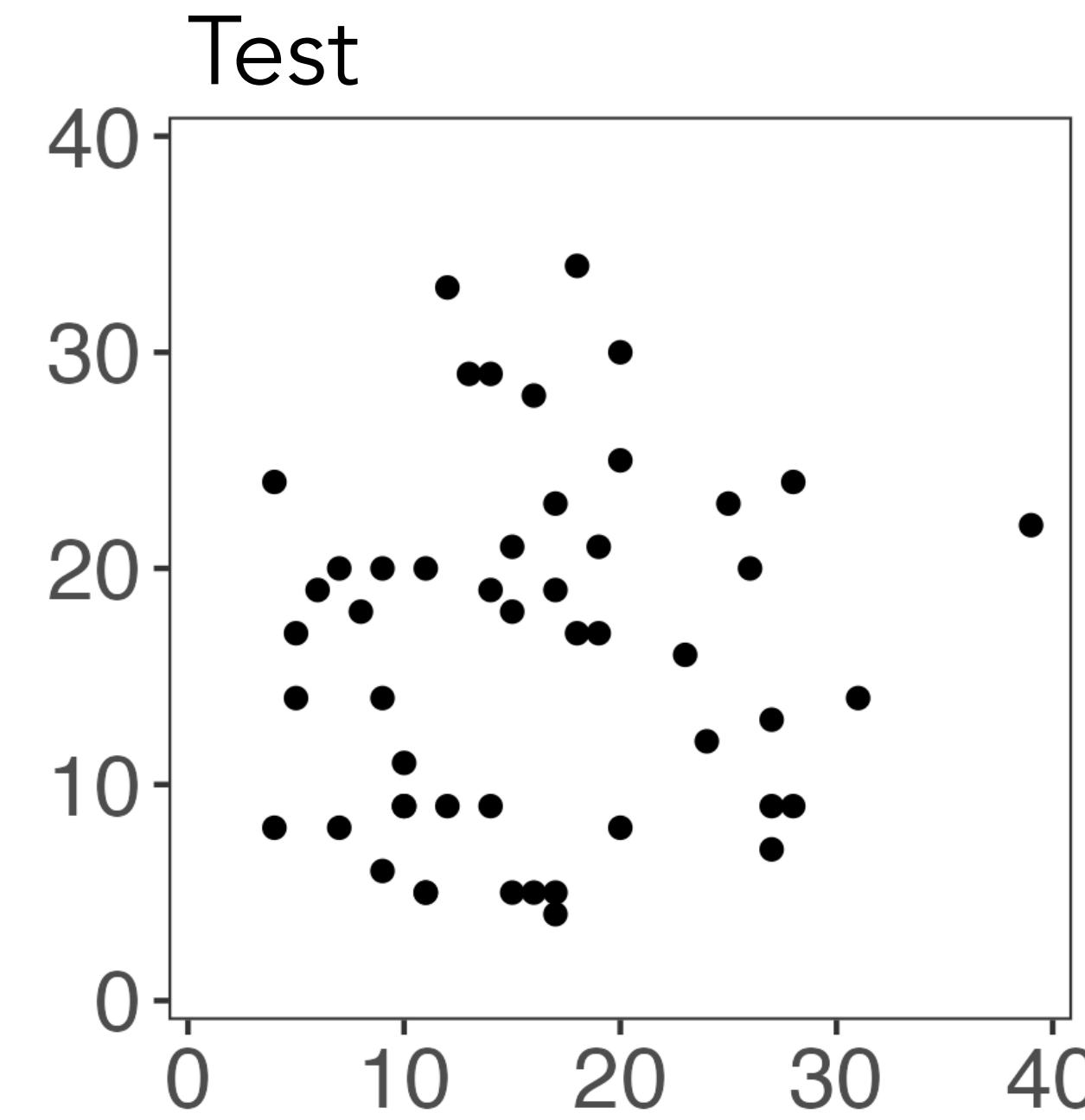
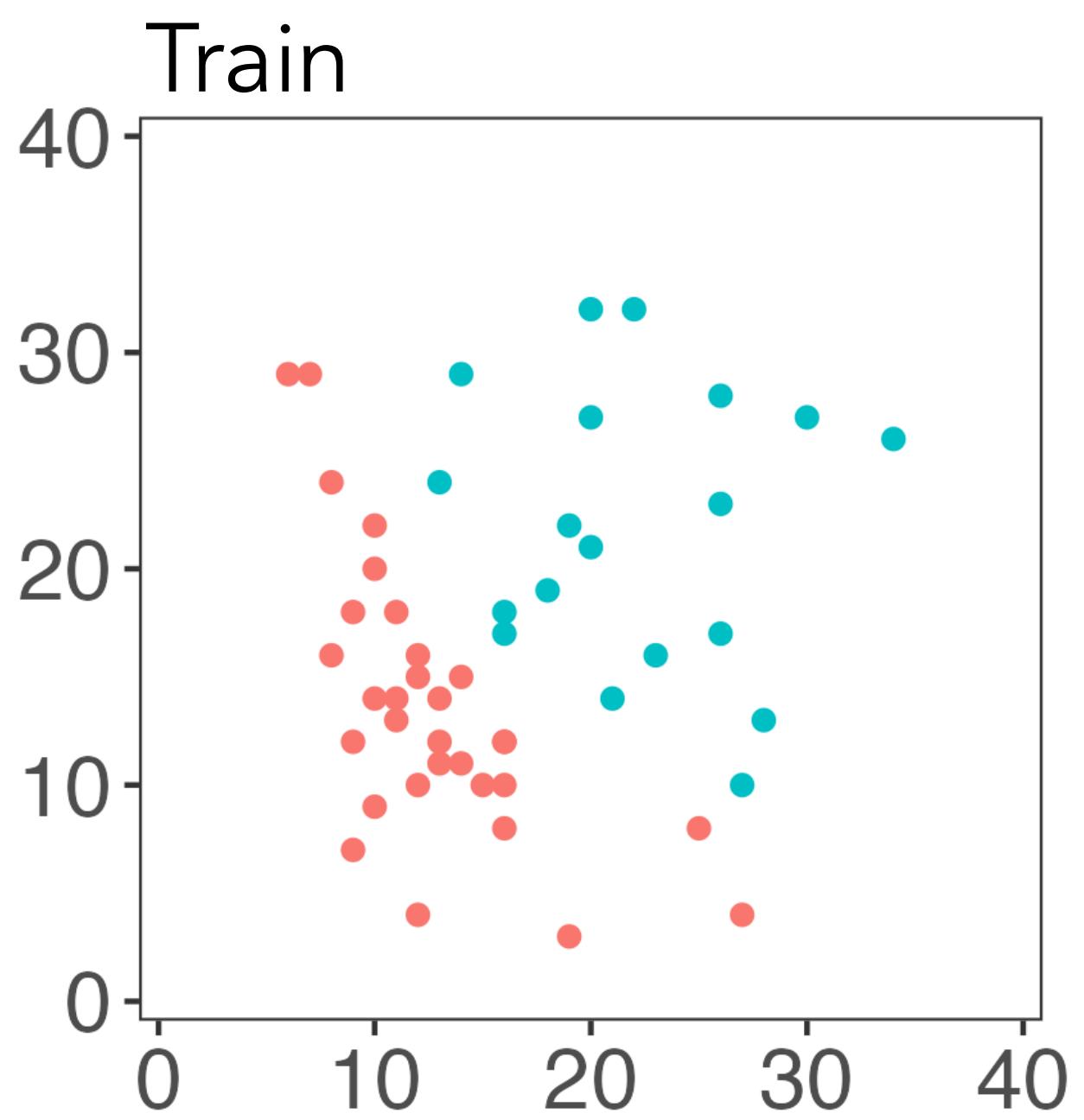
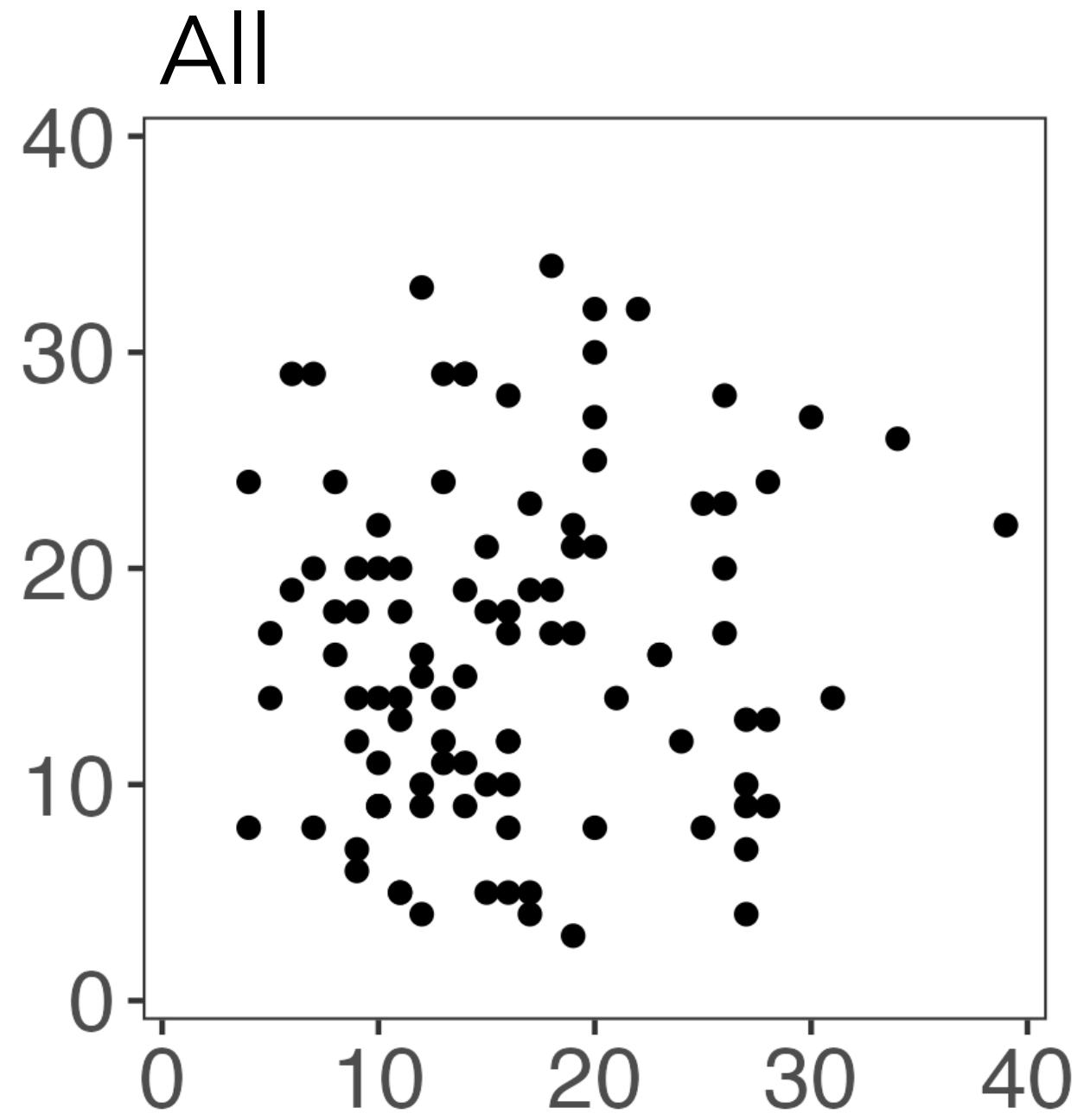


# Sample splitting cannot be used in this context



**Step 1:** split  
observations into  
train/test.

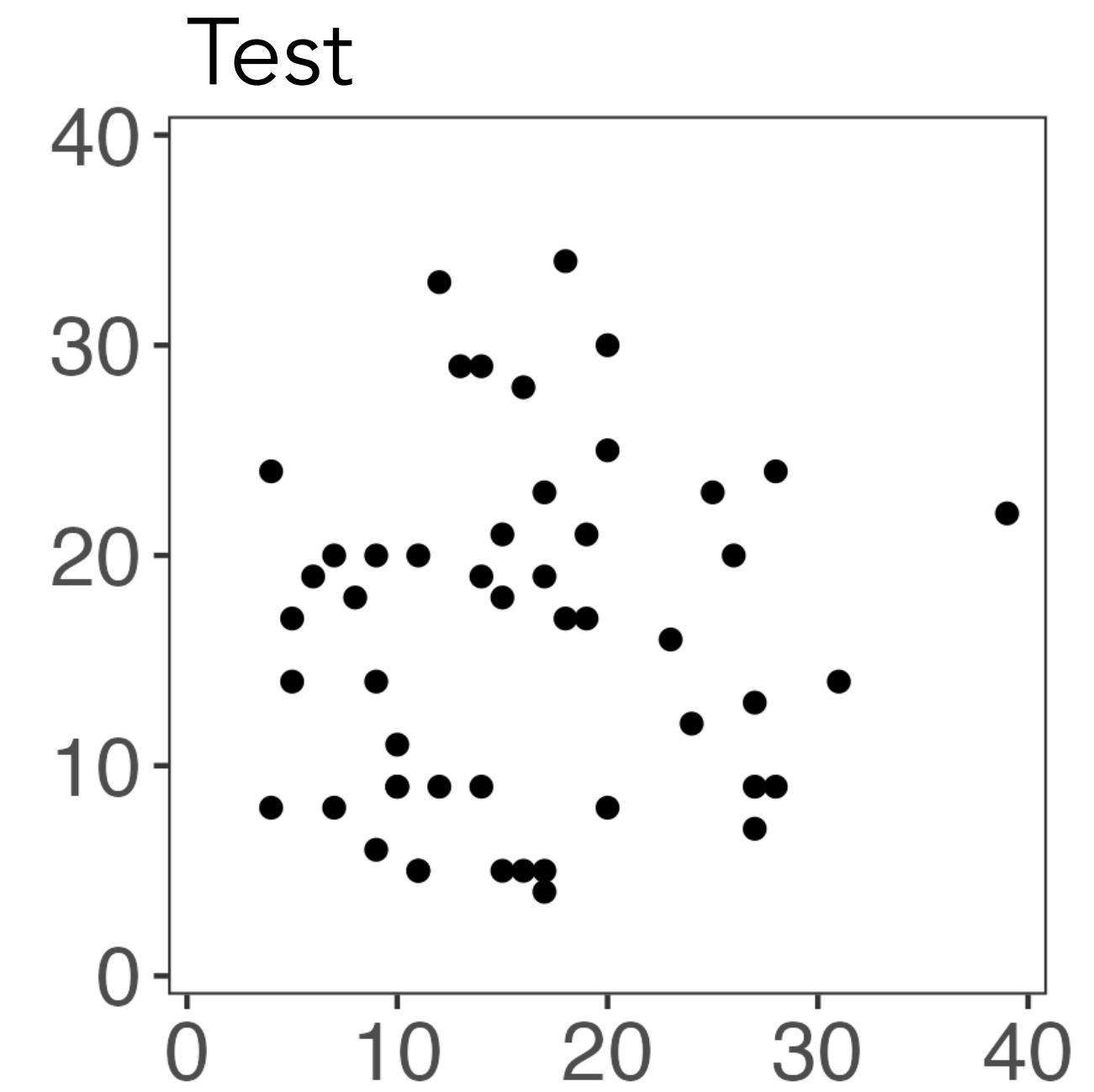
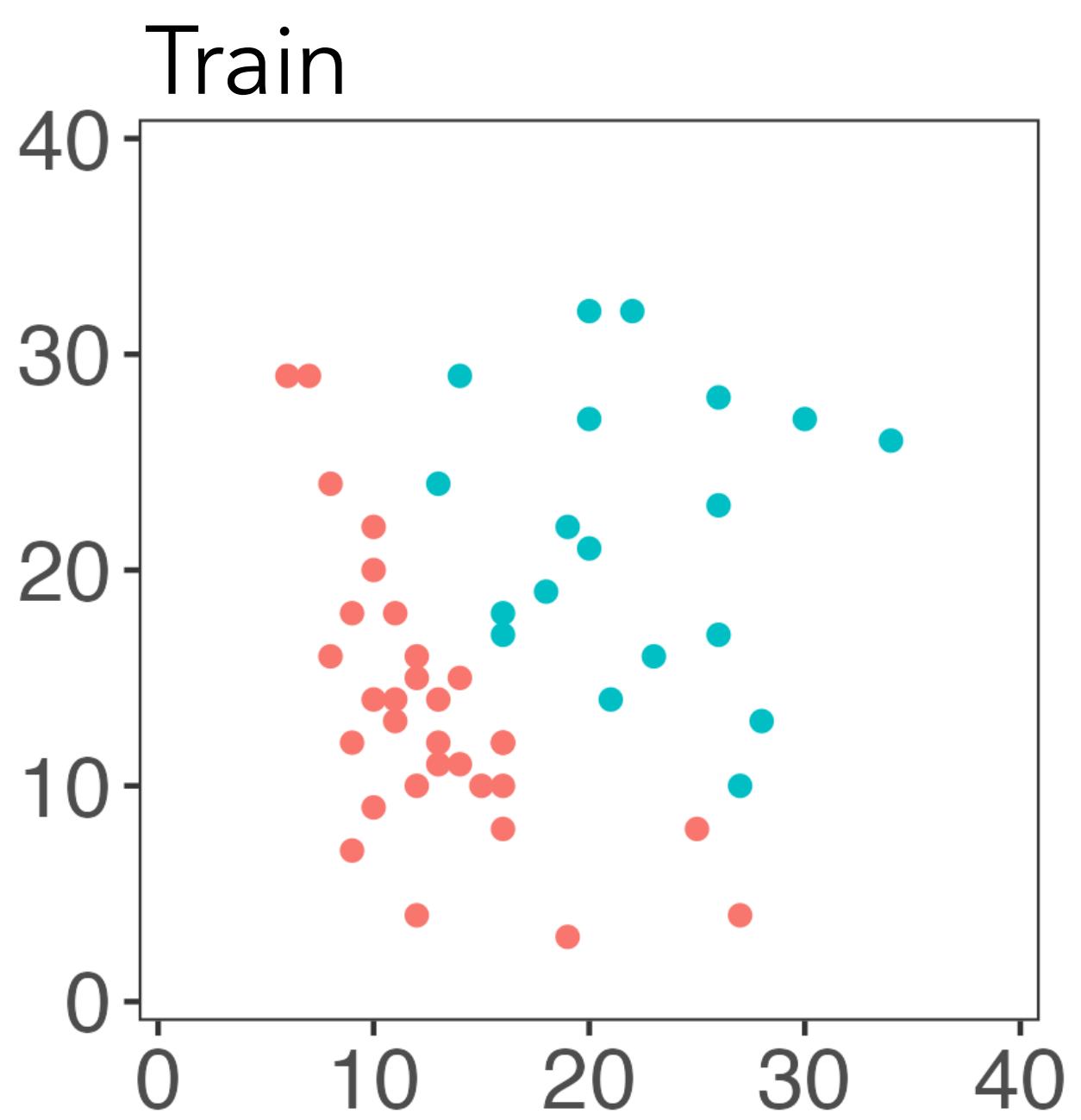
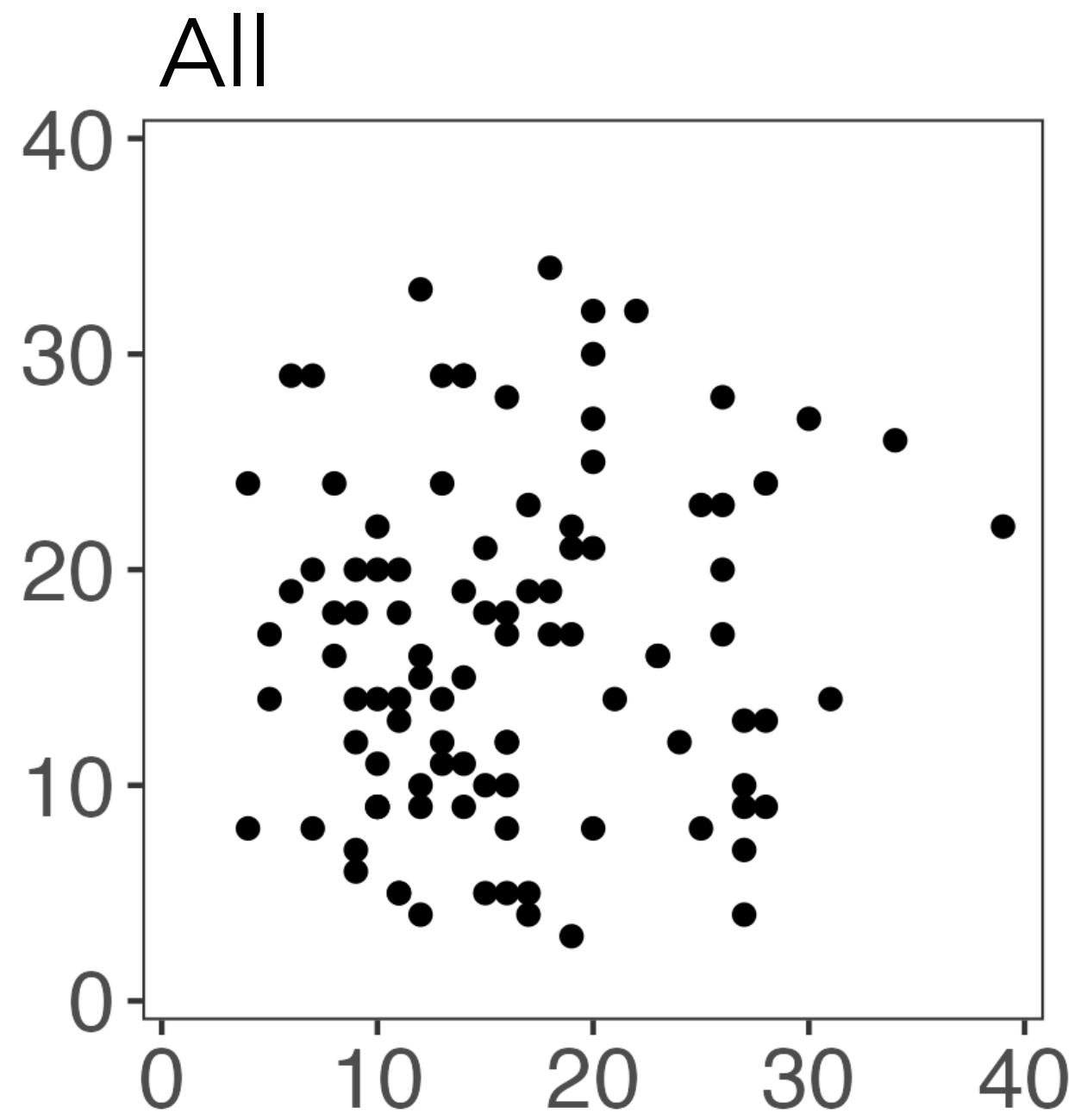
# Sample splitting cannot be used in this context



**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

# Sample splitting cannot be used in this context

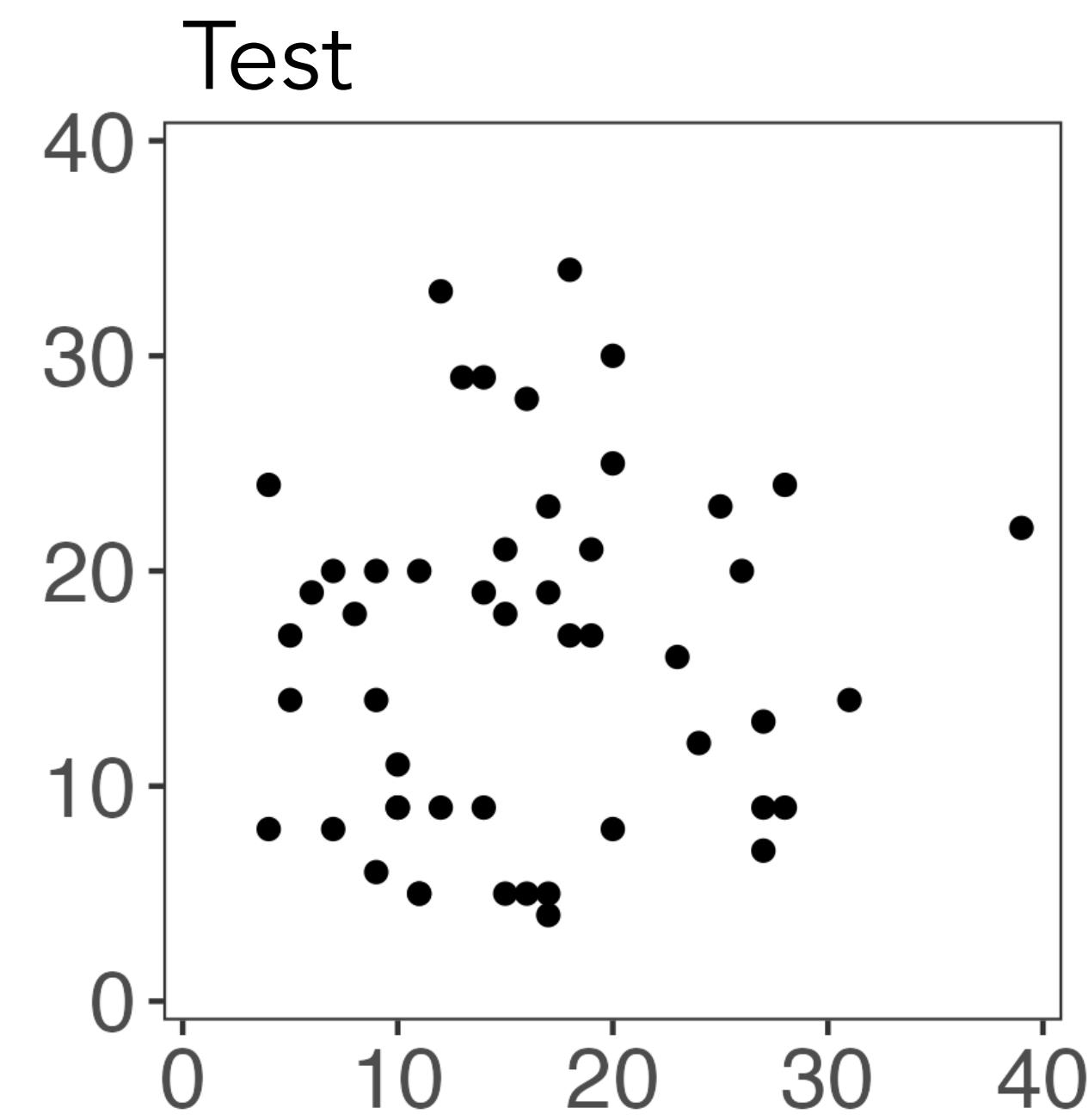
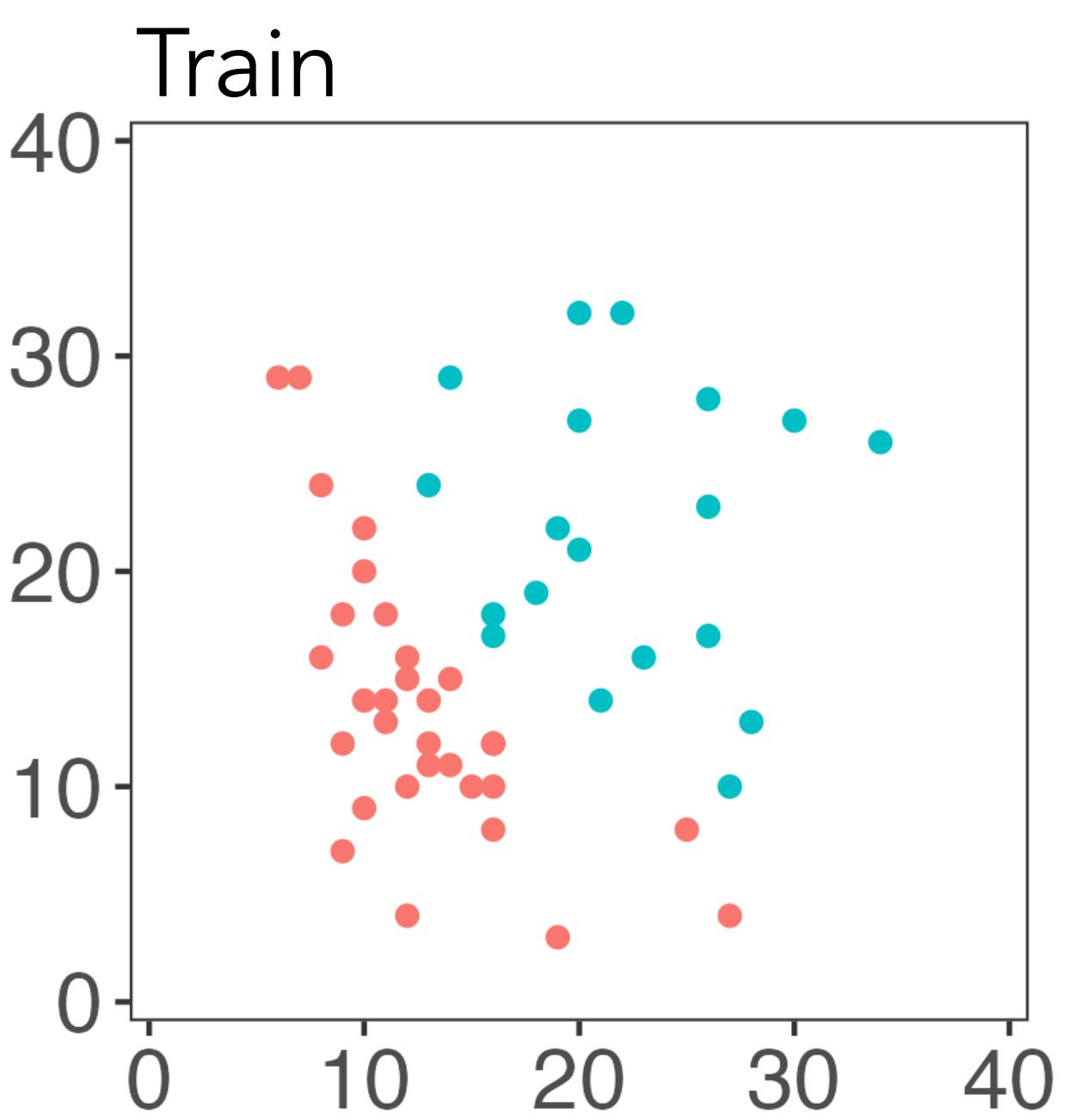
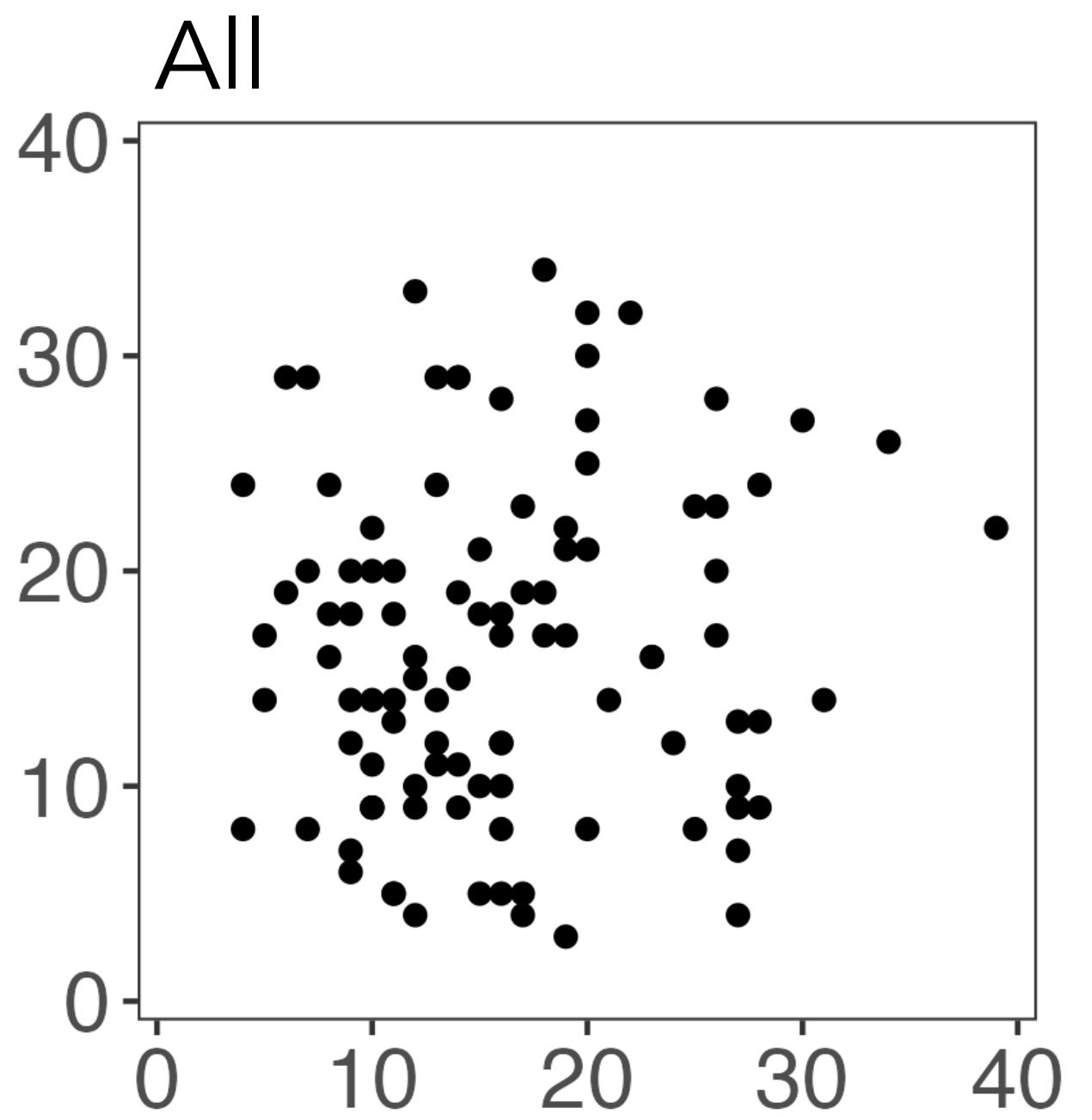


**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

**Step 3:** test for difference in means using test set.

# Sample splitting cannot be used in this context



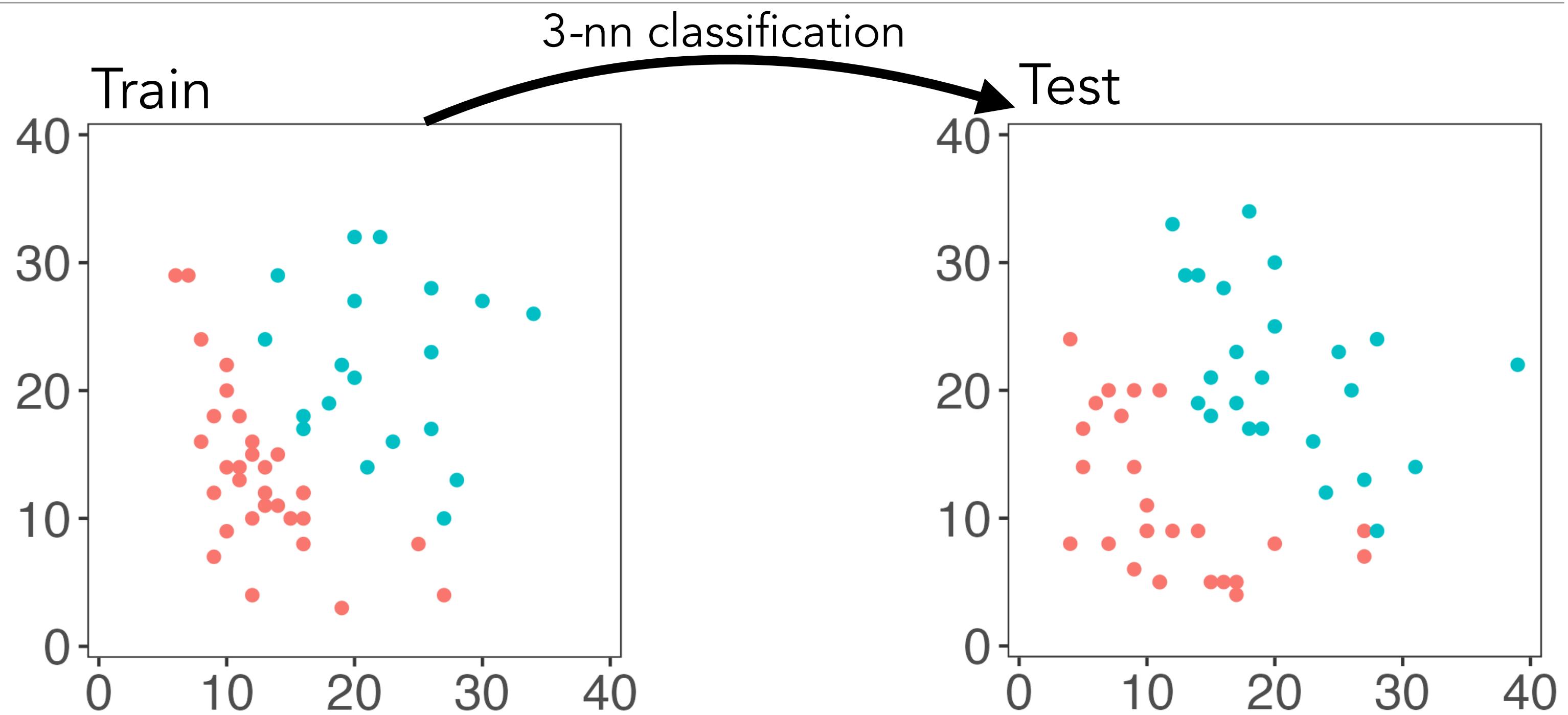
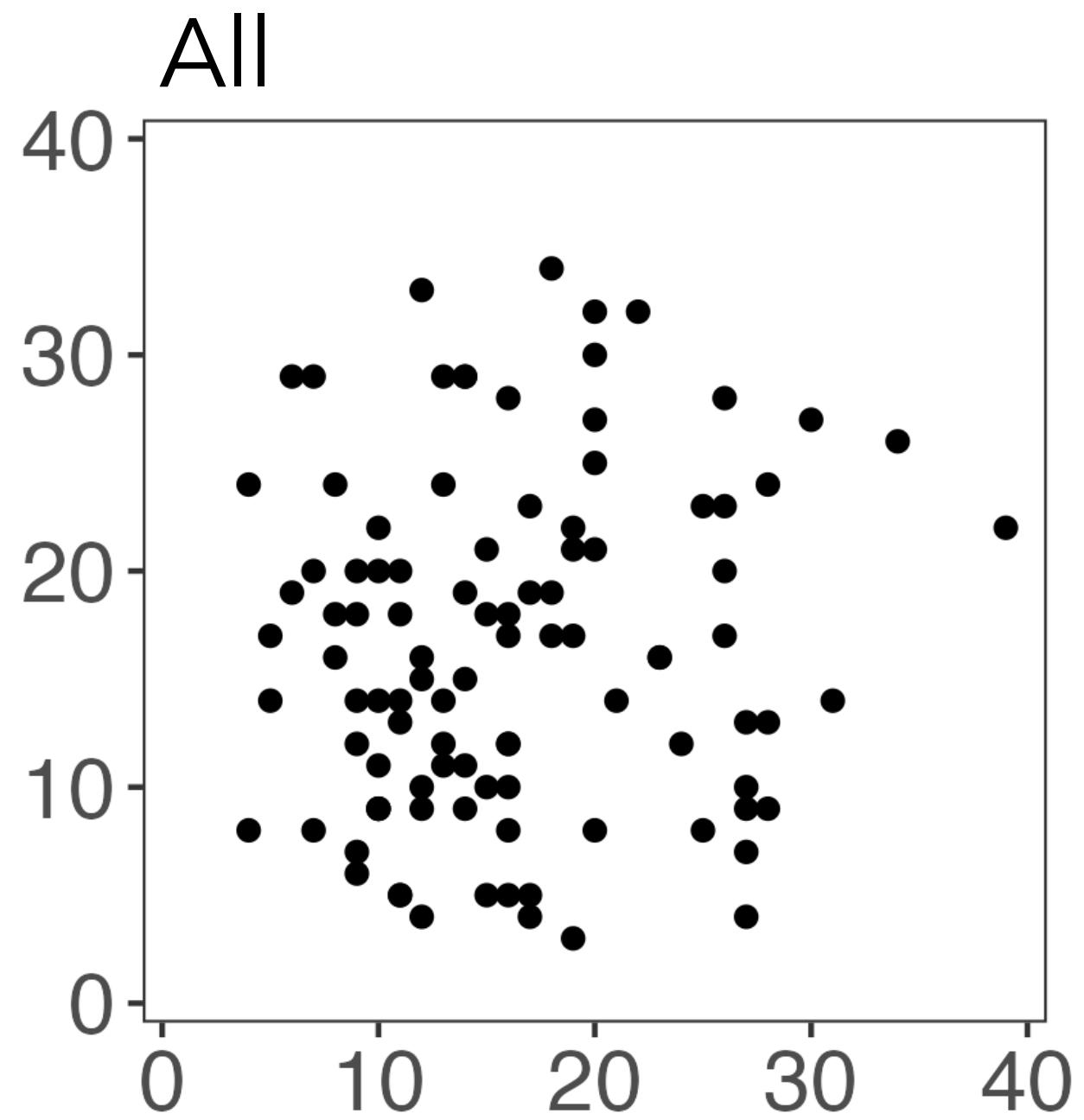
**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

**Step 2.5:** assign labels to observations in test set.

**Step 3:** test for difference in means using test set.

# Sample splitting cannot be used in this context



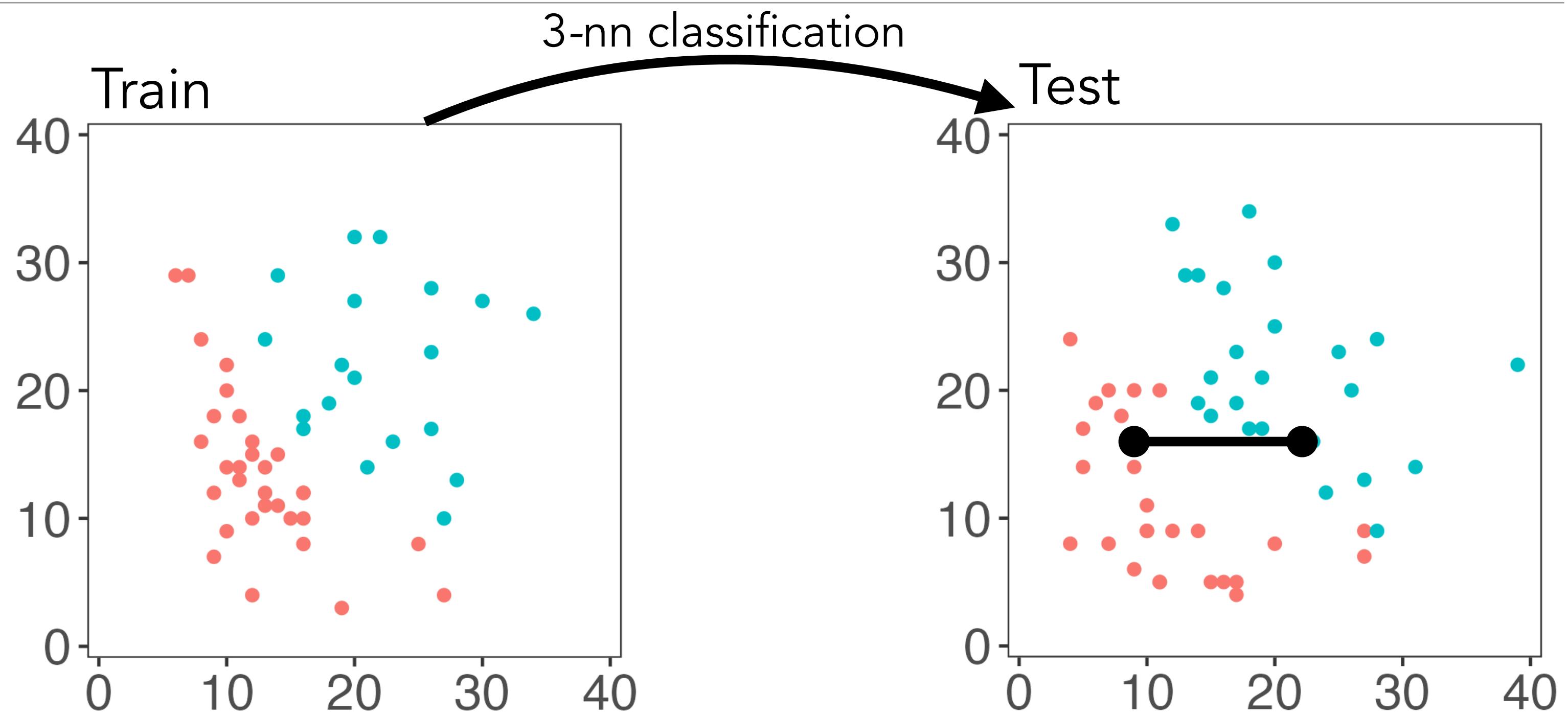
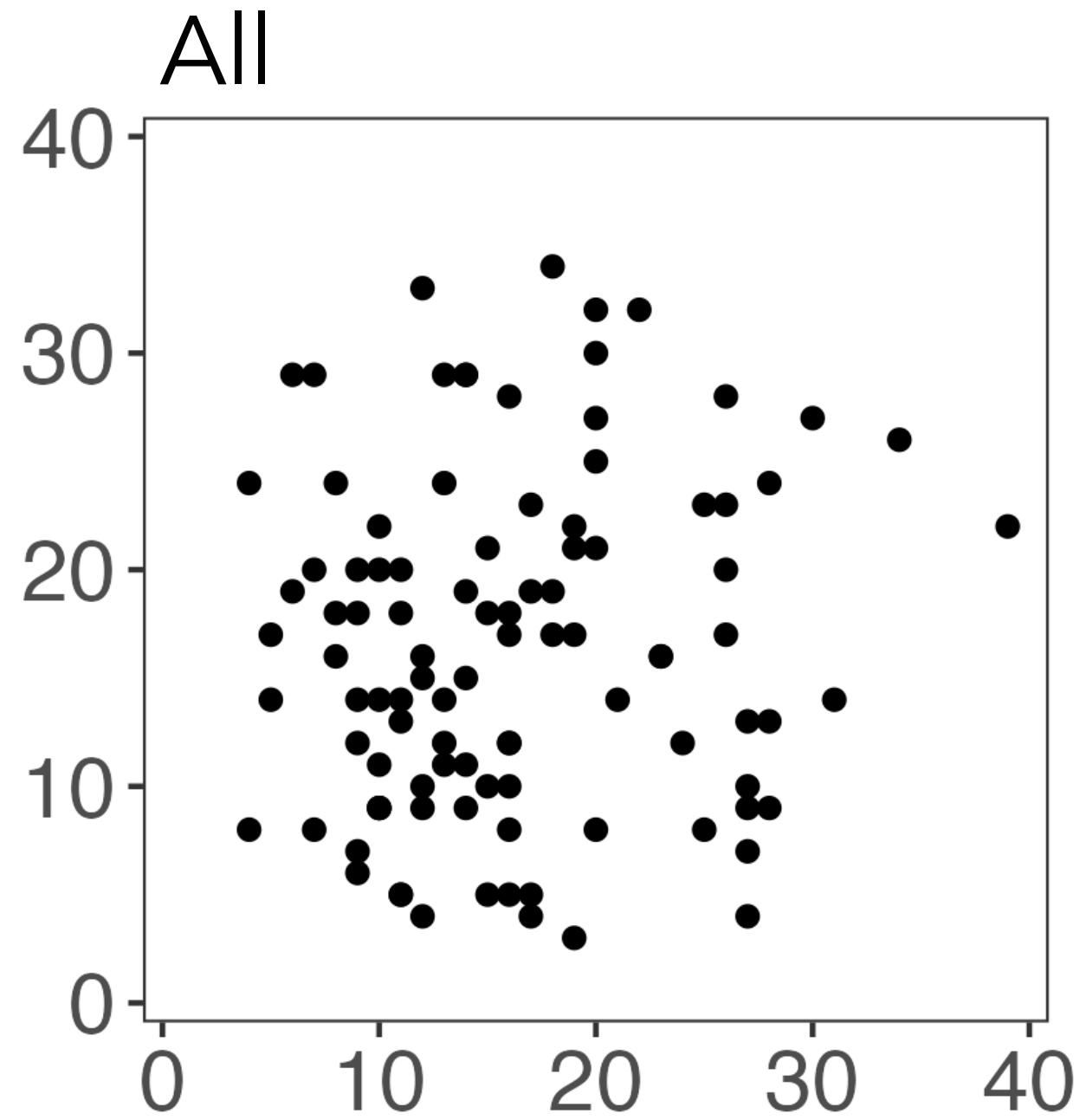
**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

**Step 2.5:** assign labels to observations in test set.

**Step 3:** test for difference in means using test set.

# Sample splitting cannot be used in this context



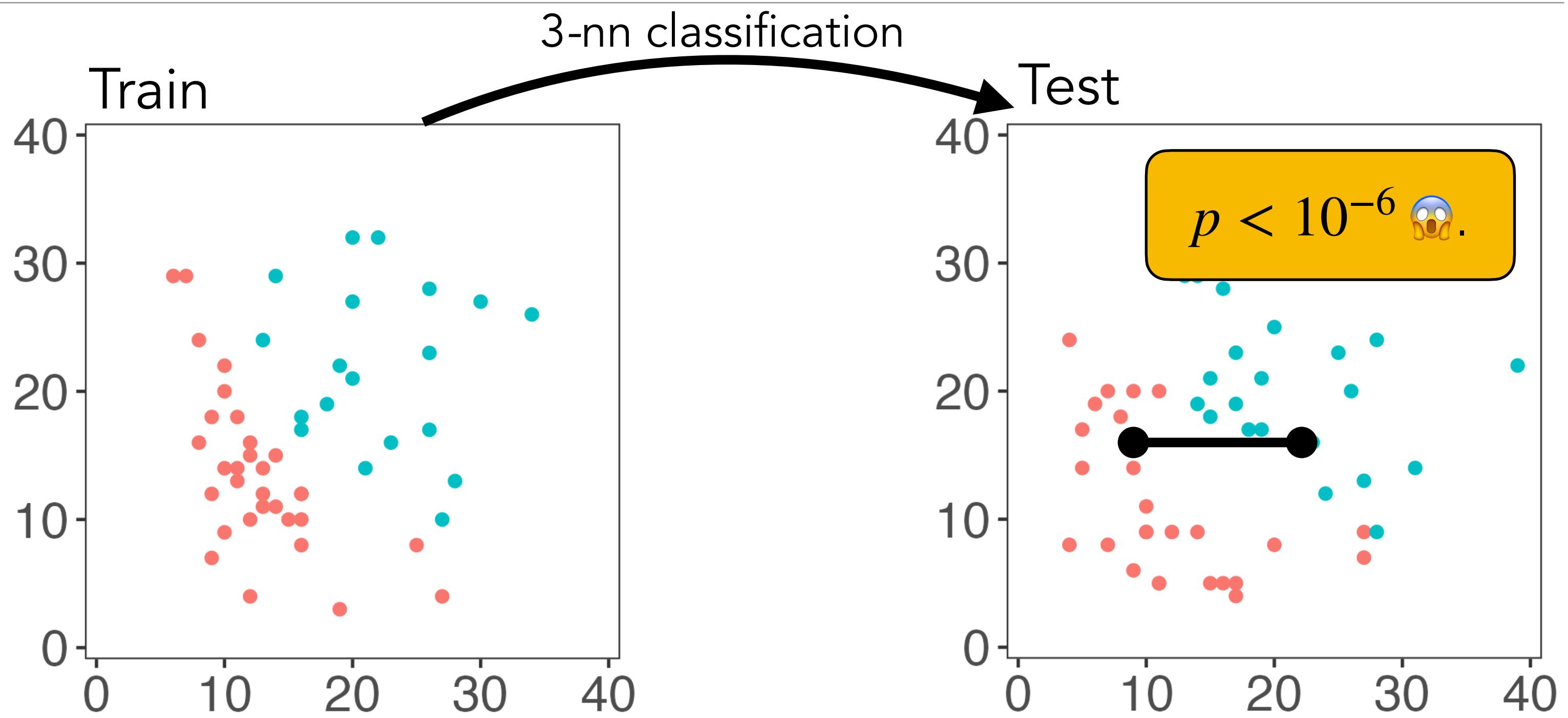
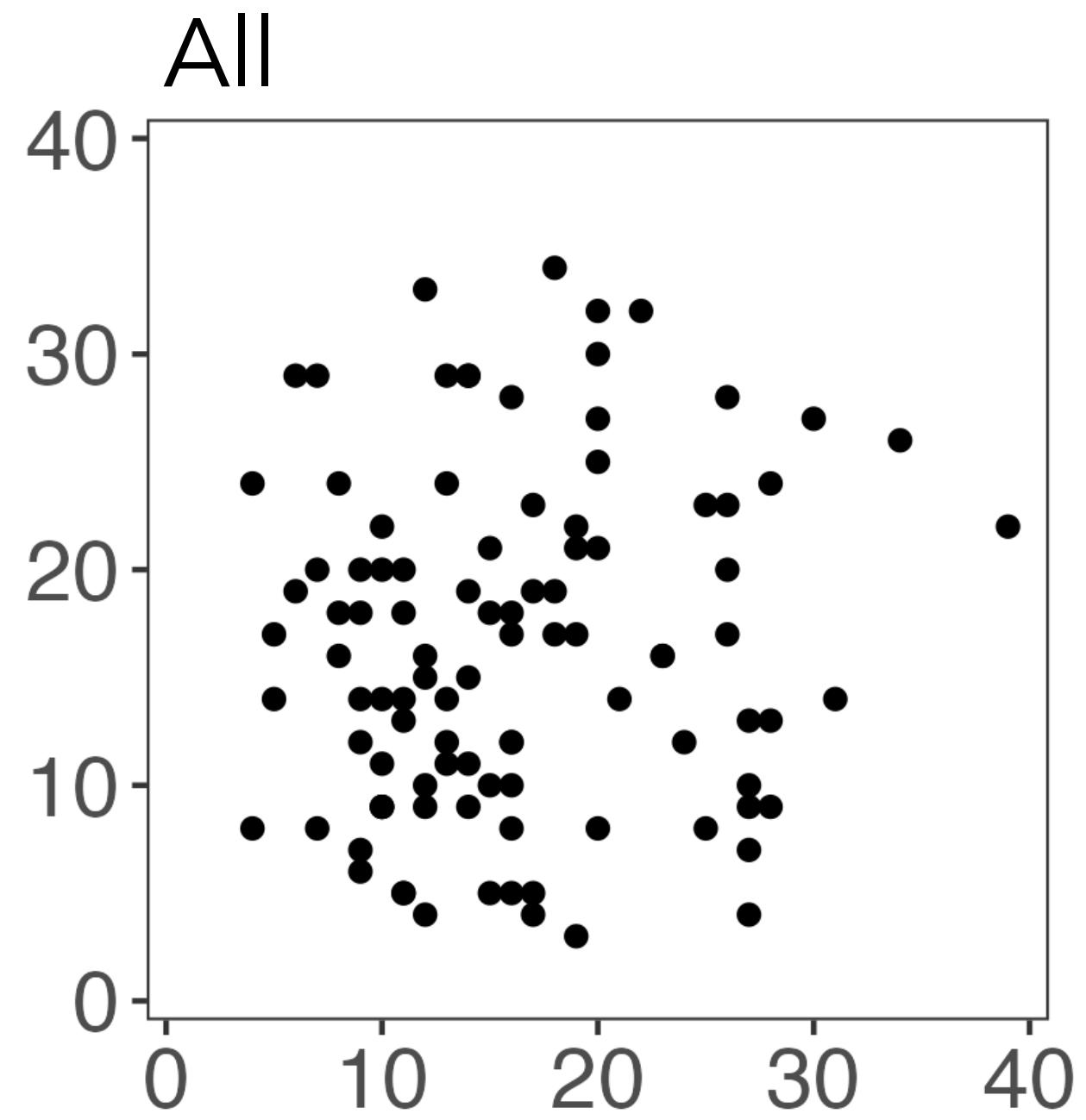
**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

**Step 2.5:** assign labels to observations in test set.

**Step 3:** test for difference in means using test set.

# Sample splitting cannot be used in this context



**Step 1:** split observations into train/test.

**Step 2:** cluster the training set.

**Step 2.5:** assign labels to observations in test set.

**Step 3:** test for difference in means using test set.

# Common scRNA-seq analysis pipelines ignore the double dipping problem

Lähnemann et al. *Genome Biology* (2020) 21:31  
<https://doi.org/10.1186/s13059-020-1926-6>

Genome Biology

REVIEW Open Access



## Eleven grand challenges in single-cell data science

David Lähnemann<sup>1,2,3</sup>, Johannes Köster<sup>1,4</sup>, Ewa Szczurek<sup>5</sup>, Davis J. McCarthy<sup>6,7</sup>, Stephanie C. Hicks<sup>8</sup>, Mark D. Robinson<sup>9</sup> , Catalina A. Vallejos<sup>10,11</sup>, Kieran R. Campbell<sup>12,13,14</sup>, Niko Beerenwinkel<sup>15,16</sup>, Ahmed Mahfouz<sup>17,18</sup>, Luca Pinello<sup>19,20,21</sup>, Pavel Skums<sup>22</sup>, Alexandros Stamatakis<sup>23,24</sup>, Camille Stephan-Otto Attolini<sup>25</sup>, Samuel Aparicio<sup>13,26</sup>, Jasmijn Baaijens<sup>27</sup>, Marleen Balvert<sup>27,28</sup>, Buys de Barbanson<sup>29,30,31</sup>, Antonio Cappuccio<sup>32</sup>, Giacomo Corleone<sup>33</sup>, Bas E. Dutilh<sup>28,34</sup>, Maria Florescu<sup>29,30,31</sup>, Victor Guryev<sup>35</sup>, Rens Holmer<sup>36</sup>, Katharina Jahn<sup>15,16</sup>, Thamar Jessurun Lobo<sup>35</sup>, Emma M. Keizer<sup>37</sup>, Indu Khatri<sup>38</sup>, Szymon M. Kielbasa<sup>39</sup>, Jan O. Korbel<sup>40</sup>, Alexey M. Kozlov<sup>23</sup>, Tzu-Hao Kuo<sup>3</sup>, Boudewijn P.F. Lelieveldt<sup>41,42</sup>, Ion I. Mandriu<sup>43</sup>, John C. Marioni<sup>44,45,46</sup>, Tobias Marschall<sup>47,48</sup>, Felix Mölder<sup>1,49</sup>, Amir Niknejad<sup>50,51</sup>, Lukasz Raczkowski<sup>5</sup>, Marcel Reinders<sup>17,18</sup>, Jeroen de Ridder<sup>29,30</sup>, Antoine-Emmanuel Saliba<sup>52</sup>, Antonios Somarakis<sup>42</sup>, Oliver Stegle<sup>40,46,53</sup>, Fabian J. Theis<sup>54</sup>, Huan Yang<sup>55</sup>, Alex Zelikovsky<sup>56,57</sup>, Alice C. McHardy<sup>3</sup>, Benjamin J. Raphael<sup>58</sup>, Sohrab P. Shah<sup>59</sup> and Alexander Schönhuth<sup>27,28\*</sup>

### Status

Currently, the vast majority of differential expression detection methods assume that the groups of cells to be compared are known in advance (e.g., experimental conditions or cell types). However, current analysis pipelines typically rely on clustering or cell type assignment to identify such groups, before downstream differential analysis is performed, without propagating the uncertainty in these assignments or accounting for the double use of data (clustering, differential testing between clusters).

# Common scRNA-seq analysis pipelines ignore the double dipping problem

Lähnemann et al. *Genome Biology* (2020) 21:31  
<https://doi.org/10.1186/s13059-020-1926-6>

Genome Biology

**REVIEW** **Open Access**



## Eleven grand challenges in single-cell data science

David Lähnemann<sup>1,2,3</sup>, Johannes Köster<sup>1,4</sup>, Ewa Szczurek<sup>5</sup>, Davis J. McCarthy<sup>6,7</sup>, Stephanie C. Hicks<sup>8</sup>, Mark D. Robinson<sup>9</sup> , Catalina A. Vallejos<sup>10,11</sup>, Kieran R. Campbell<sup>12,13,14</sup>, Niko Beerenwinkel<sup>15,16</sup>, Ahmed Mahfouz<sup>17,18</sup>, Luca Pinello<sup>19,20,21</sup>, Pavel Skums<sup>22</sup>, Alexandros Stamatakis<sup>23,24</sup>, Camille Stephan-Otto Attolini<sup>25</sup>, Samuel Aparicio<sup>13,26</sup>, Jasmijn Baaijens<sup>27</sup>, Marleen Balvert<sup>27,28</sup>, Buys de Barbanson<sup>29,30,31</sup>, Antonio Cappuccio<sup>32</sup>, Giacomo Corleone<sup>33</sup>, Bas E. Dutilh<sup>28,34</sup>, Maria Florescu<sup>29,30,31</sup>, Victor Guryev<sup>35</sup>, Rens Holmer<sup>36</sup>, Katharina Jahn<sup>15,16</sup>, Thamar Jessurun Lobo<sup>35</sup>, Emma M. Keizer<sup>37</sup>, Indu Khatri<sup>38</sup>, Szymon M. Kielbasa<sup>39</sup>, Jan O. Korbel<sup>40</sup>, Alexey M. Kozlov<sup>23</sup>, Tzu-Hao Kuo<sup>3</sup>, Boudewijn P.F. Lelieveldt<sup>41,42</sup>, Ion I. Mandriu<sup>43</sup>, John C. Marioni<sup>44,45,46</sup>, Tobias Marschall<sup>47,48</sup>, Felix Mölder<sup>1,49</sup>, Amir Niknejad<sup>50,51</sup>, Lukasz Raczkowski<sup>5</sup>, Marcel Reinders<sup>17,18</sup>, Jeroen de Ridder<sup>29,30</sup>, Antoine-Emmanuel Saliba<sup>52</sup>, Antonios Somarakis<sup>42</sup>, Oliver Stegle<sup>40,46,53</sup>, Fabian J. Theis<sup>54</sup>, Huan Yang<sup>55</sup>, Alex Zelikovsky<sup>56,57</sup>, Alice C. McHardy<sup>3</sup>, Benjamin J. Raphael<sup>58</sup>, Sohrab P. Shah<sup>59</sup> and Alexander Schönhuth<sup>27,28\*</sup>

**Status**

Currently, the vast majority of differential expression detection methods assume that the groups of cells to be compared are known in advance (e.g., experimental conditions or cell types). However, current analysis pipelines typically rely on clustering or cell type assignment to identify such groups, before downstream differential analysis is performed, without propagating the uncertainty in these assignments or accounting for the double use of data (clustering, differential testing between clusters).

Seurat **4.0.6** [Install](#) [Get started](#) [Vignettes](#) [Extensions](#) [FAQ](#) [News](#) [Reference](#) [Archive](#)

## Gene expression markers of identity classes

Source: [R/generics.R](#), [R/differential\\_expression.R](#)

Finds markers (differentially expressed genes) for identity classes

`FindMarkers(object, ...)`

## Details

p-value adjustment is performed using bonferroni correction based on the total number of genes in the dataset. Other correction methods are not recommended, as Seurat pre-filters genes using the arguments above, reducing the number of tests performed. Lastly, as Aaron Lun has pointed out, p-values should be interpreted cautiously, as the genes used for clustering are the same genes tested for differential expression.

# Outline

---

1. Motivation: settings where sample splitting doesn't work
2. **Poisson thinning**
3. Data thinning
4. Application to single-cell RNA sequencing data
5. Ongoing work

Reminder: sample splitting does not help us with our motivating example

scRNA-seq dataset

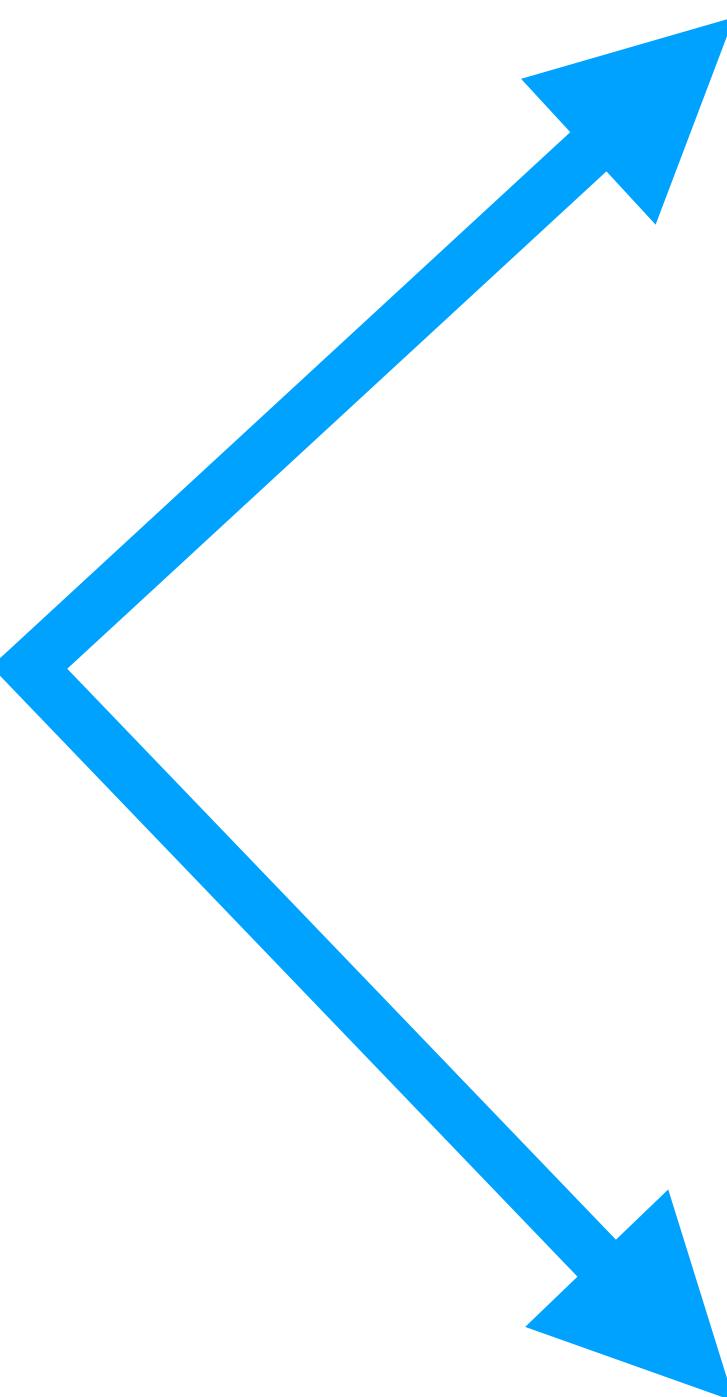
	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

Train

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	28

Test

	Gene 1	Gene 2
Cell 3	11	5
Cell 4	22	21



Reminder: sample splitting does not help us with our motivating example

scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

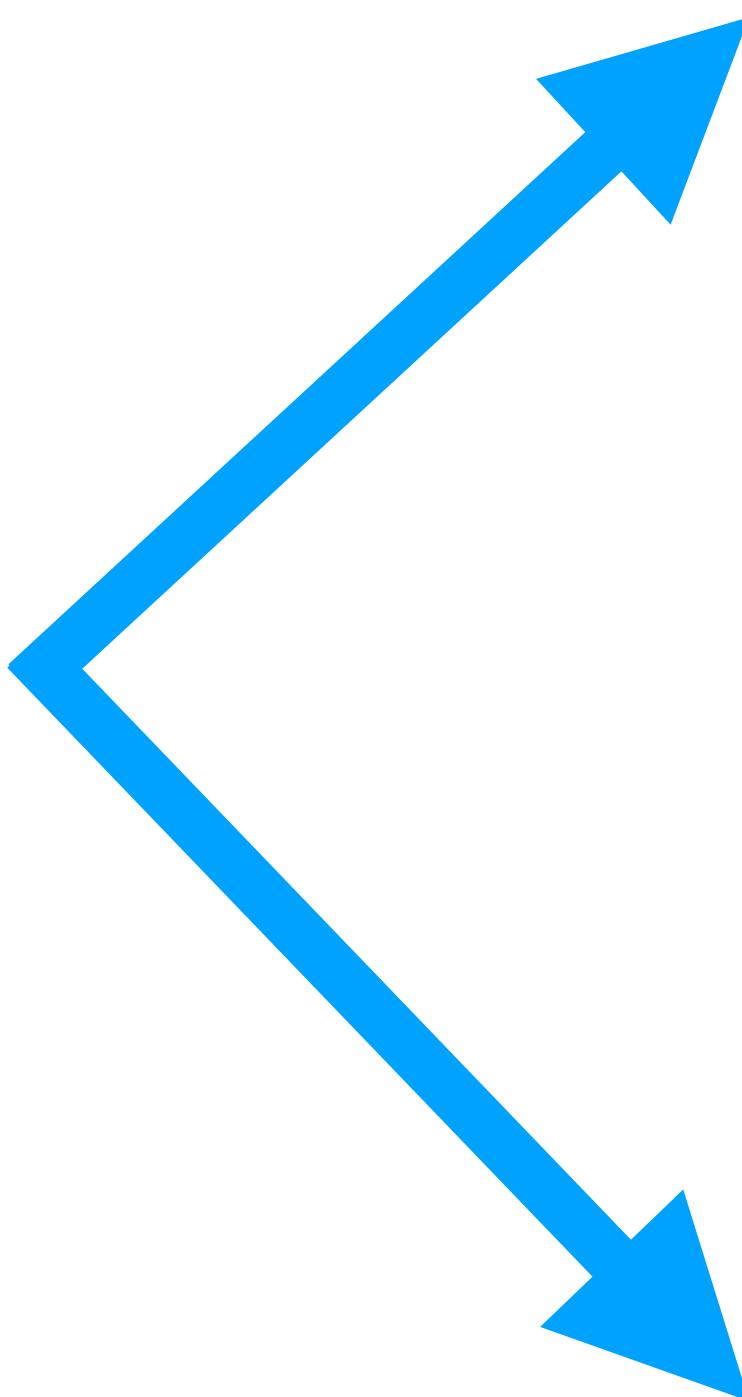
Train

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	28

Estimating clusters on training set

Test

	Gene 1	Gene 2
Cell 3	11	5
Cell 4	22	21



Reminder: sample splitting does not help us with our motivating example

scRNA-seq dataset

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

Train

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	28

Test

	Gene 1	Gene 2
Cell 3	11	5
Cell 4	22	21

Estimating clusters on training set

does not yield cluster assignments for test set.

## An alternative: Poisson thinning

---

$X$

	<b>Gene 1</b>	<b>Gene 2</b>
<b>Cell 1</b>	18	6
<b>Cell 2</b>	31	8
<b>Cell 3</b>	11	31
<b>Cell 4</b>	22	34

## An alternative: Poisson thinning

---

$X$

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

$X^{(1)}$

	Gene 1	Gene 2
Cell 1	14	1
Cell 2	10	6
Cell 3	5	17
Cell 4	6	25

$X^{(2)}$

	Gene 3	Gene 4
Cell 1	4	5
Cell 2	21	2
Cell 3	6	14
Cell 4	16	9

## An alternative: Poisson thinning

$X$

	Gene 1	Gene 2
Cell 1	18	6
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

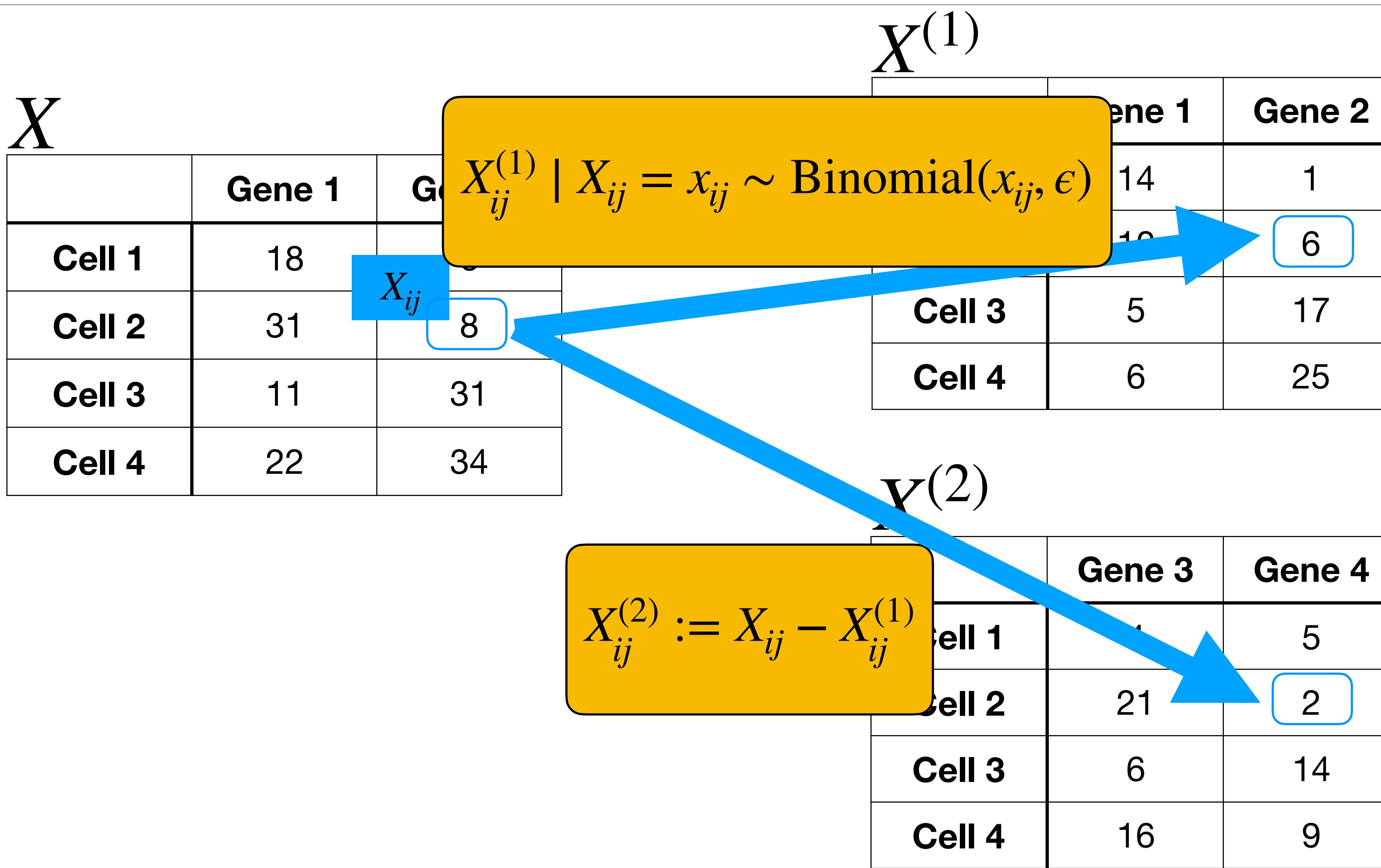
$X^{(1)}$

	Gene 1	Gene 2
Cell 1	14	1
Cell 2	10	6
Cell 3	5	17
Cell 4	6	25

$X^{(2)}$

	Gene 3	Gene 4
Cell 1	4	5
Cell 2	21	2
Cell 3	6	14
Cell 4	16	9

## An alternative: Poisson thinning



## An alternative: Poisson thinning

$X$

	Gene 1	Gene 2
Cell 1	18	$X_{ij}$
Cell 2	31	8
Cell 3	11	31
Cell 4	22	34

$X^{(1)}$

$$X_{ij}^{(1)} \mid X_{ij} = x_{ij} \sim \text{Binomial}(x_{ij}, \epsilon)$$

$X^{(2)}$

If  $X_{ij} \sim \text{Poisson}(\Lambda_{ij})$ , then:

1.  $X_{ij}^{(1)} \sim \text{Poisson}(\epsilon \Lambda_{ij})$
2.  $X_{ij}^{(2)} \sim \text{Poisson}((1 - \epsilon) \Lambda_{ij})$
3.  $X_{ij}^{(1)} \perp\!\!\!\perp X_{ij}^{(2)}$

Gene 1  
Gene 2

14 1

10 6

5 17

6 25

Gene 3  
Gene 4

4 5

21 2

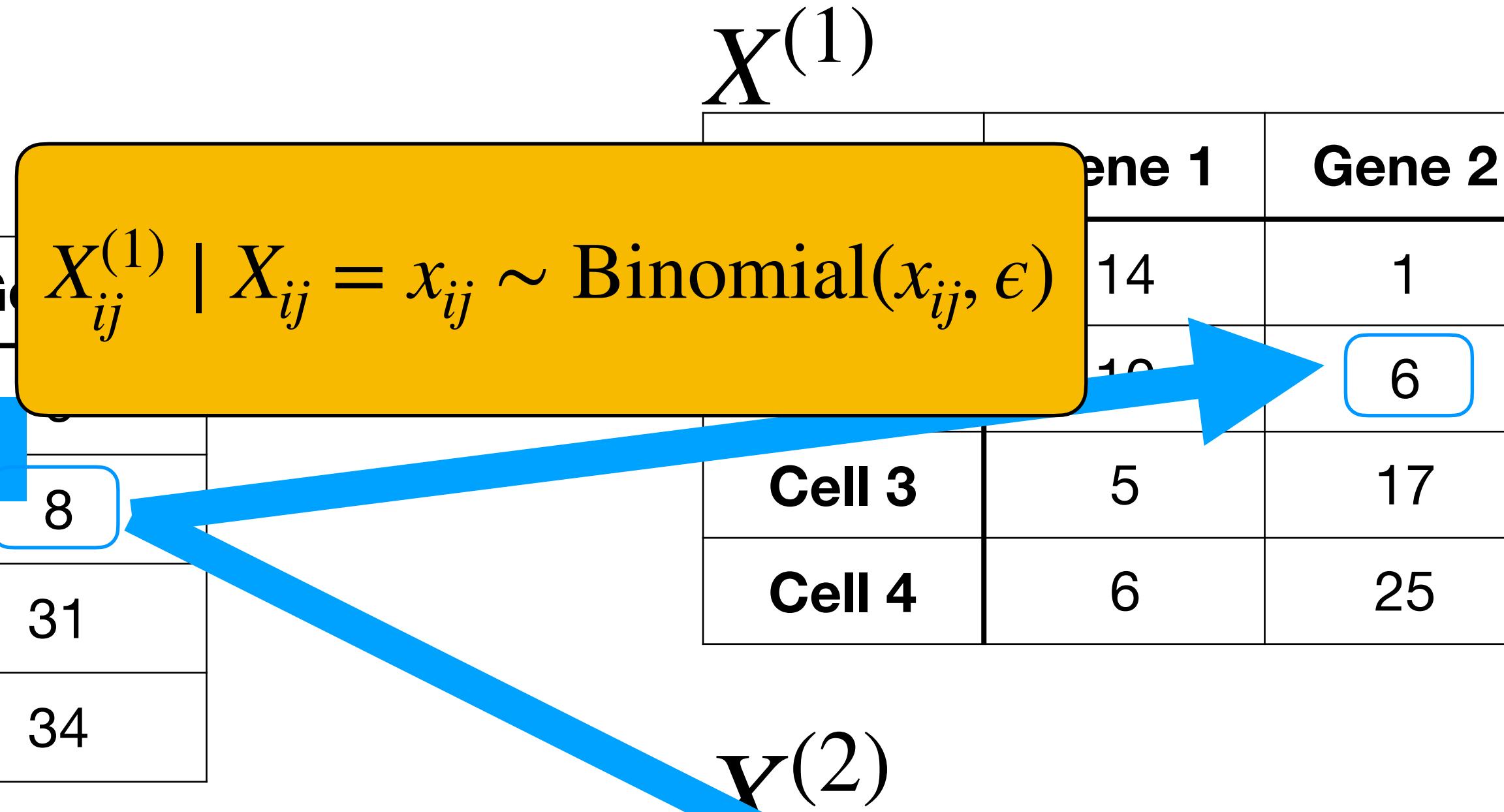
6 14

16 9

A very well-known result.

## An alternative: Poisson thinning

	<b>Gene 1</b>	<b>Gene 2</b>
<b>Cell 1</b>	18	14
<b>Cell 2</b>	31	1
<b>Cell 3</b>	11	5
<b>Cell 4</b>	22	17
	31	25



Estimate clusters.

If  $X_{ij} \sim \text{Poisson}(\Lambda_{ij})$ , then:

1.  $X_{ij}^{(1)} \sim \text{Poisson}(\epsilon \Lambda_{ij})$
2.  $X_{ij}^{(2)} \sim \text{Poisson}((1 - \epsilon) \Lambda_{ij})$
3.  $X_{ij}^{(1)} \perp\!\!\!\perp X_{ij}^{(2)}$

$X^{(2)}$

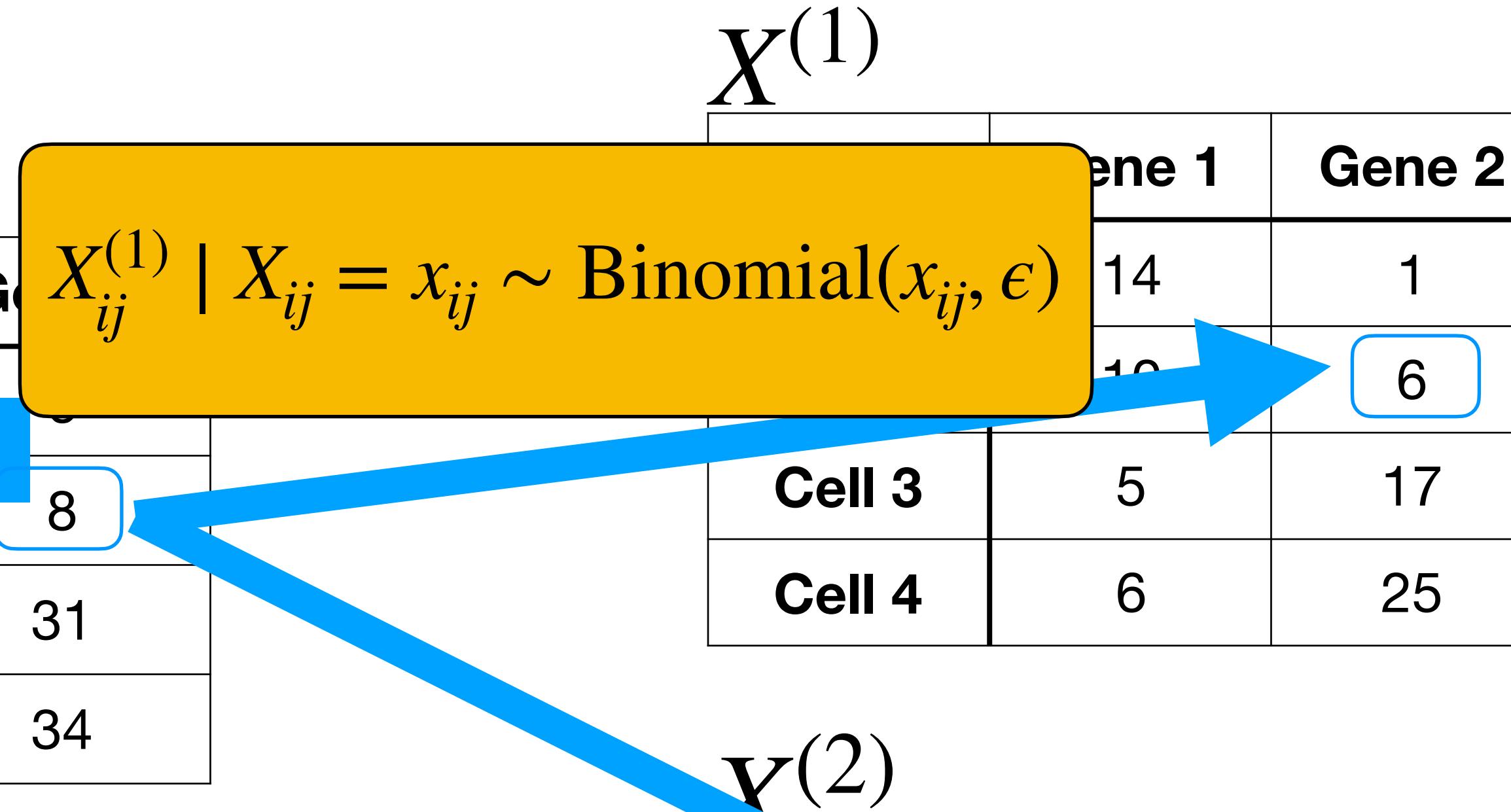
	<b>Gene 3</b>	<b>Gene 4</b>
<b>Cell 1</b>	1	5
<b>Cell 2</b>	21	2
<b>Cell 3</b>	6	14
<b>Cell 4</b>	16	9

$X_{ij}^{(2)} := X_{ij} - X_{ij}^{(1)}$

A very well-known result.

# An alternative: Poisson thinning

	<b>Gene 1</b>	<b>Gene 2</b>
<b>Cell 1</b>	18	14
<b>Cell 2</b>	31	1
<b>Cell 3</b>	11	5
<b>Cell 4</b>	22	17
	31	25



Estimate clusters.

If  $X_{ij} \sim \text{Poisson}(\Lambda_{ij})$ , then:

1.  $X_{ij}^{(1)} \sim \text{Poisson}(\epsilon \Lambda_{ij})$
2.  $X_{ij}^{(2)} \sim \text{Poisson}((1 - \epsilon) \Lambda_{ij})$
3.  $X_{ij}^{(1)} \perp\!\!\!\perp X_{ij}^{(2)}$

$X^{(2)}$

	<b>Gene 3</b>	<b>Gene 4</b>
<b>Cell 1</b>	4	5
<b>Cell 2</b>	21	2
<b>Cell 3</b>	6	14
<b>Cell 4</b>	16	9
	14	13

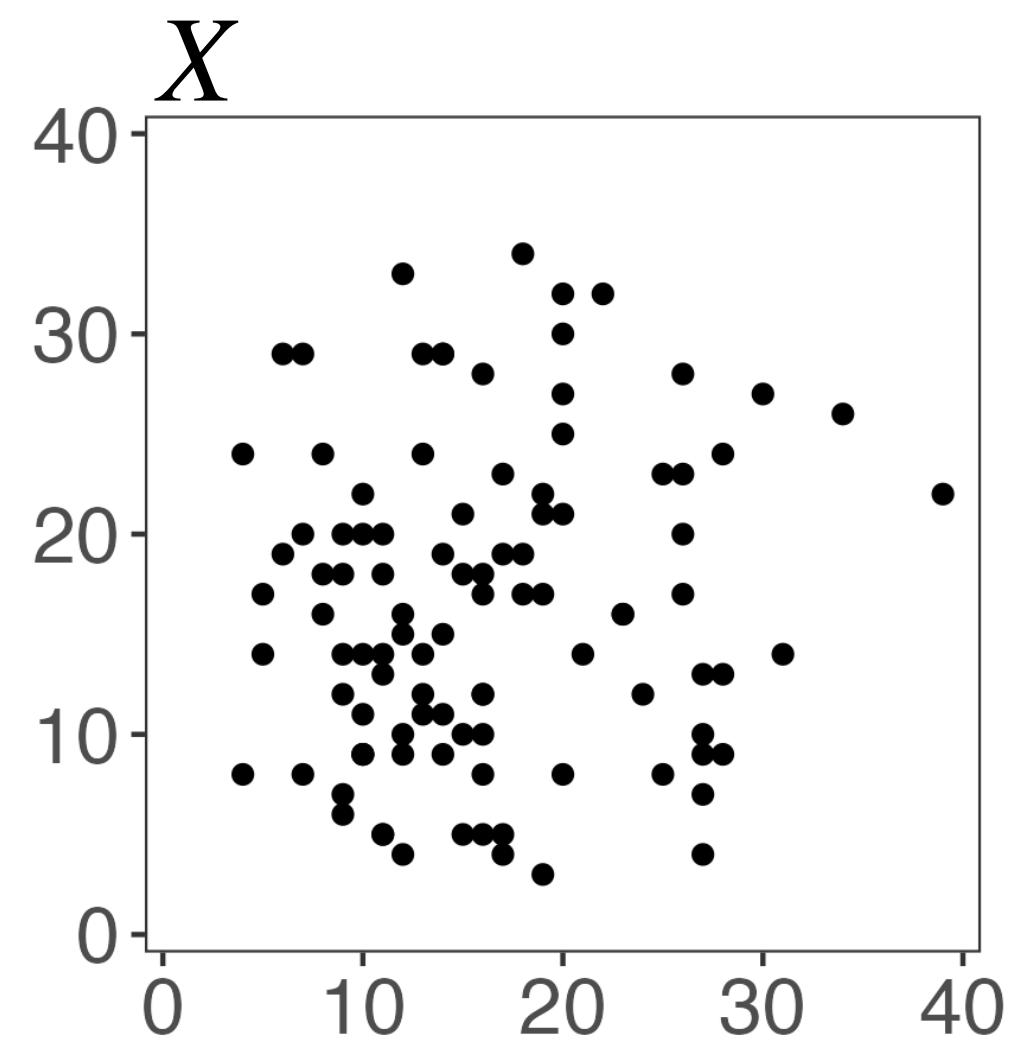
$X_{ij}^{(2)} := X_{ij} - X_{ij}^{(1)}$

Test for  
difference in  
means.

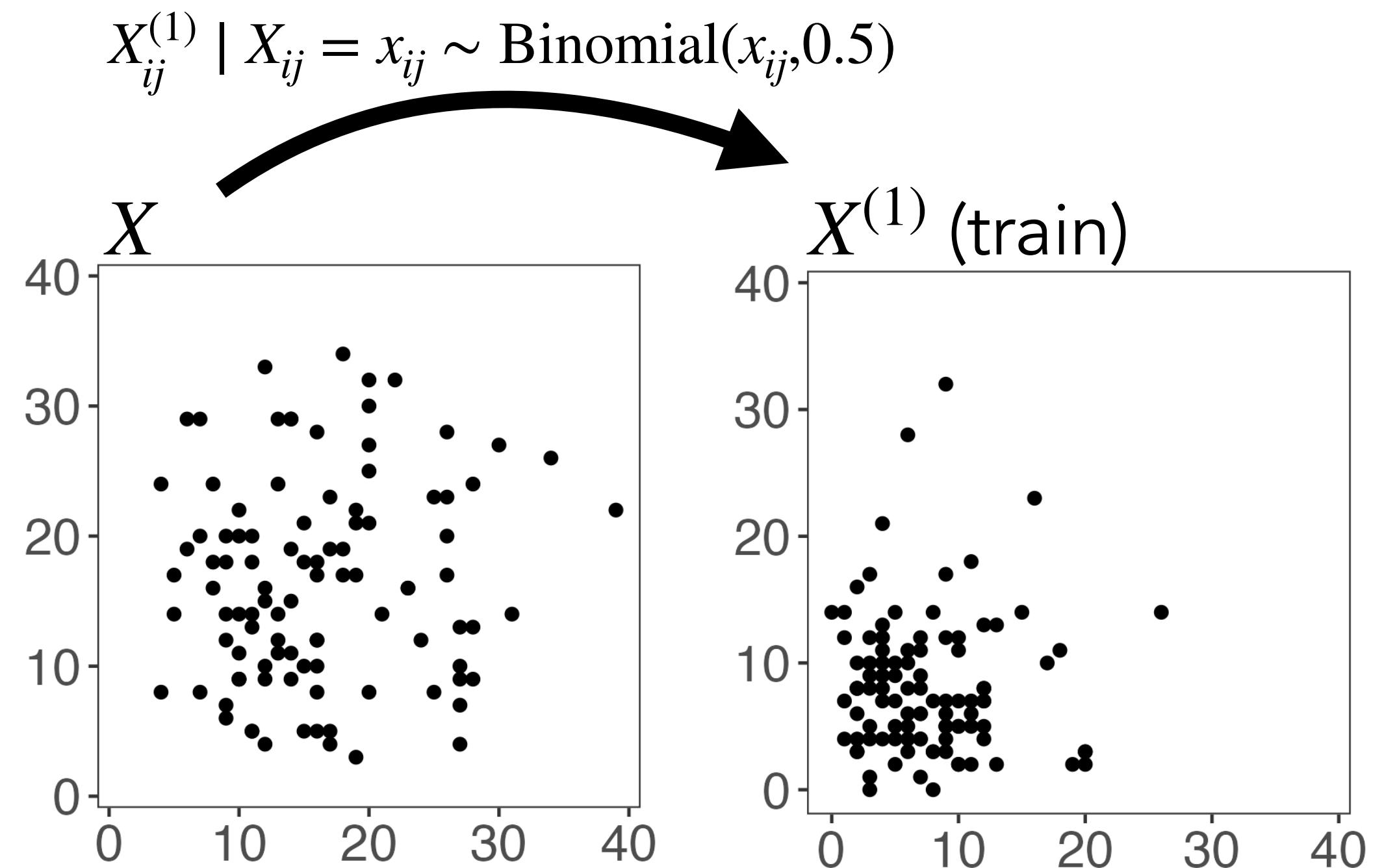
A very well-known result.

Thinning avoids the pitfall of sample splitting on our motivating examples

---

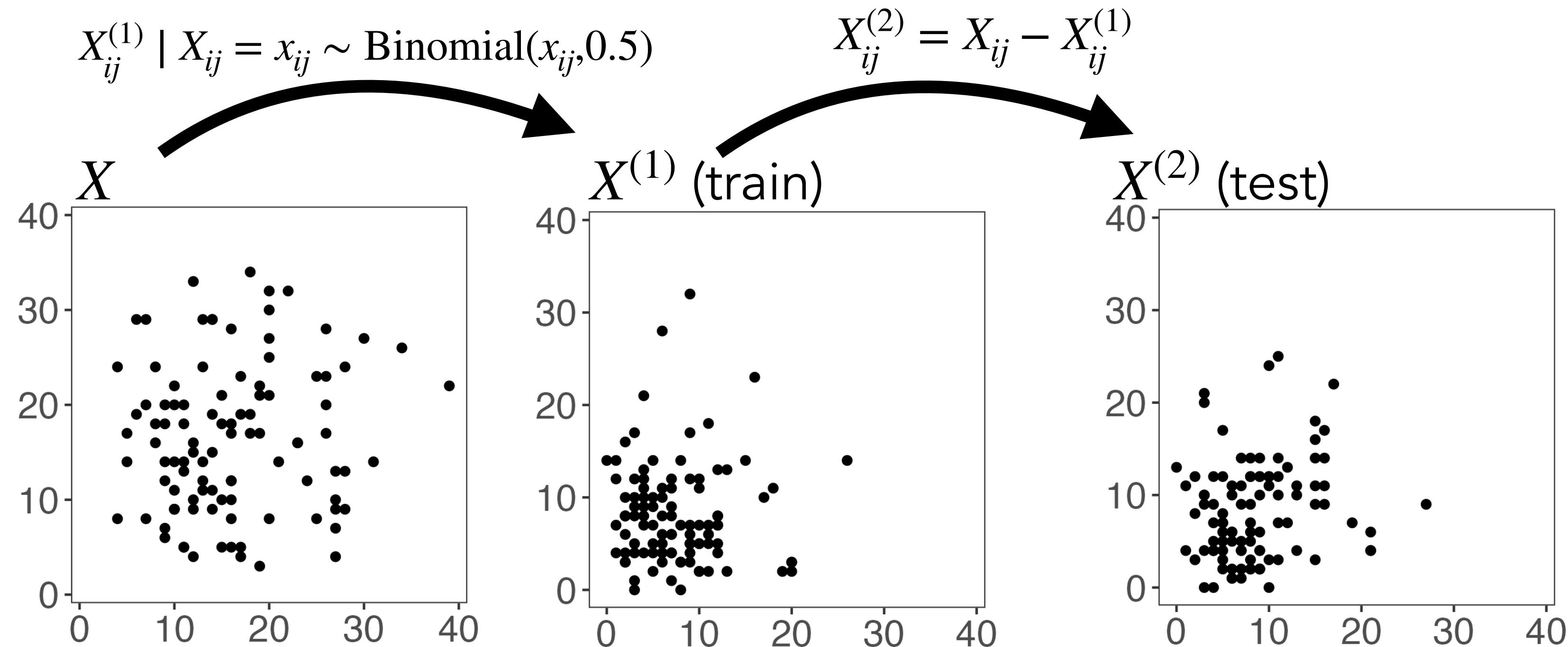


# Thinning avoids the pitfall of sample splitting on our motivating examples



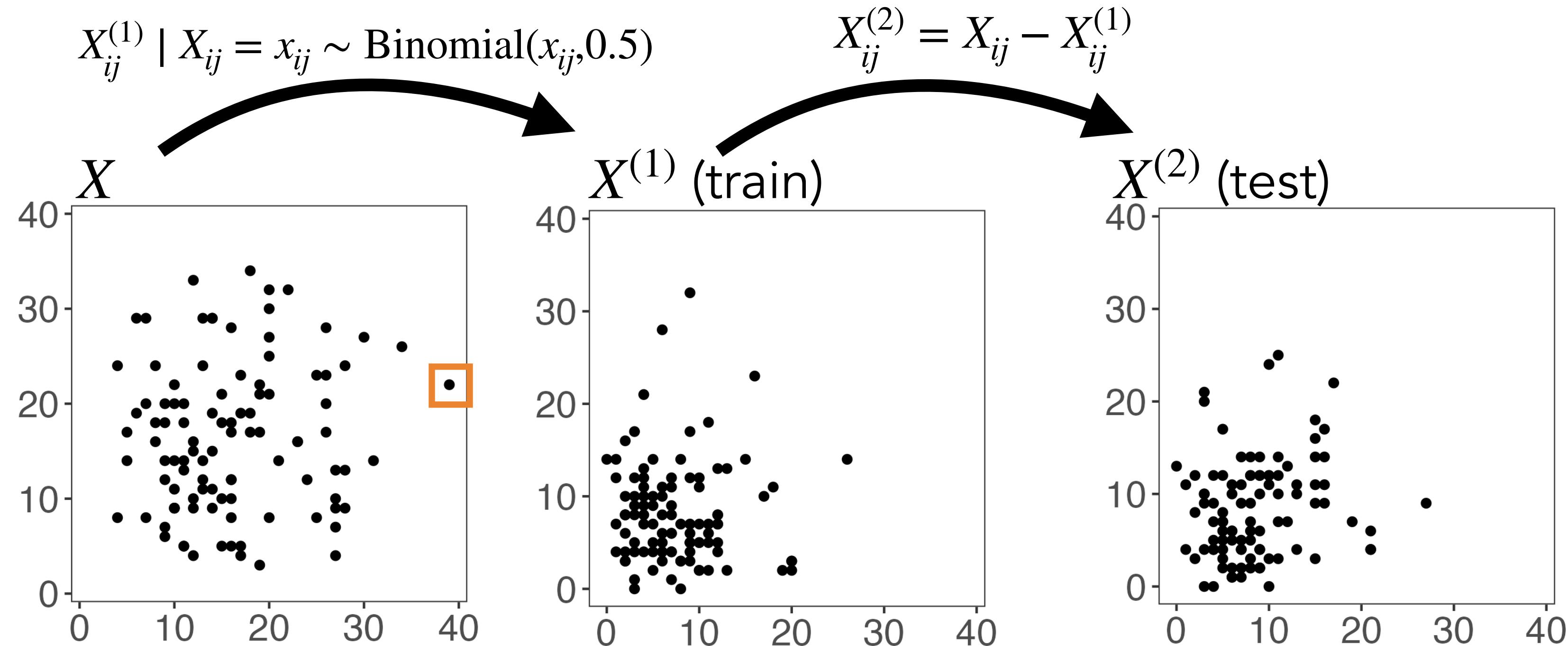
**Step 1:** thin observations into train/test.

# Thinning avoids the pitfall of sample splitting on our motivating examples



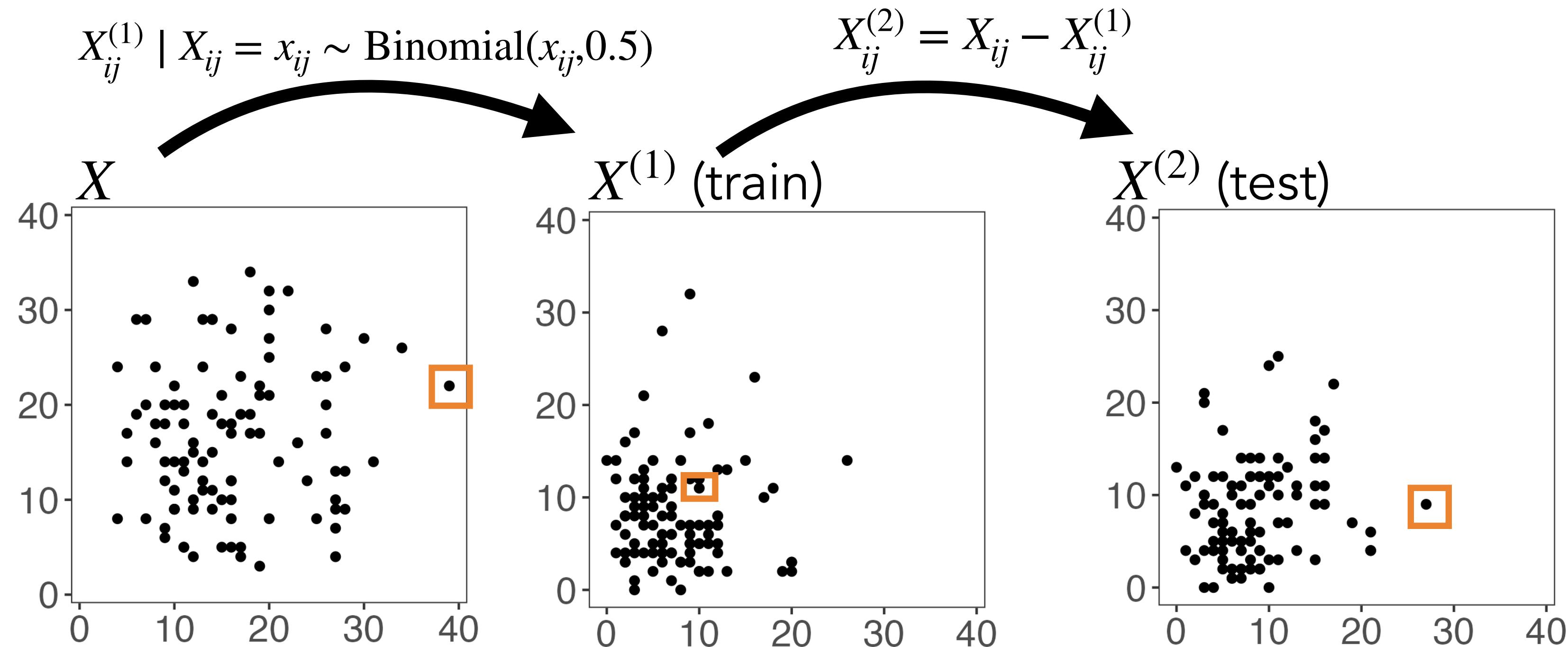
**Step 1:** thin observations into train/test.

# Thinning avoids the pitfall of sample splitting on our motivating examples



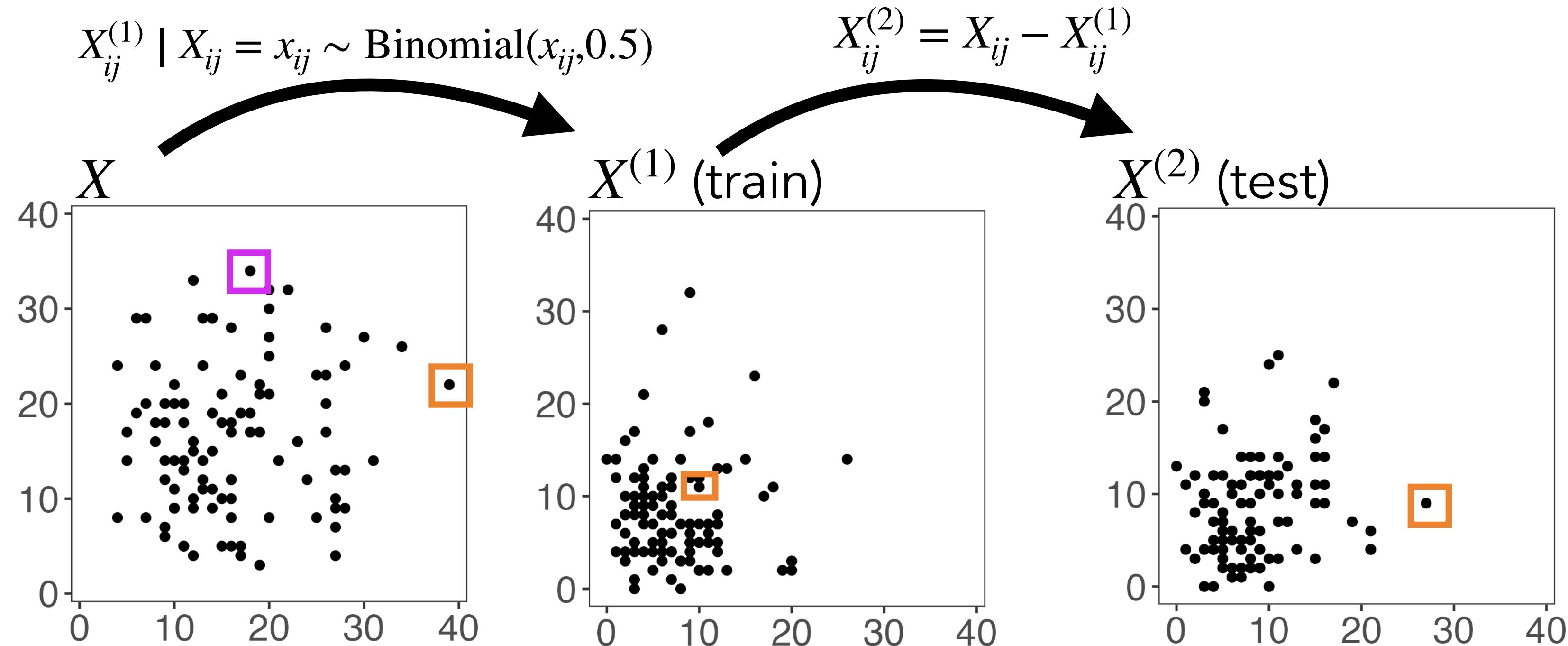
**Step 1:** thin observations into train/test.

# Thinning avoids the pitfall of sample splitting on our motivating examples



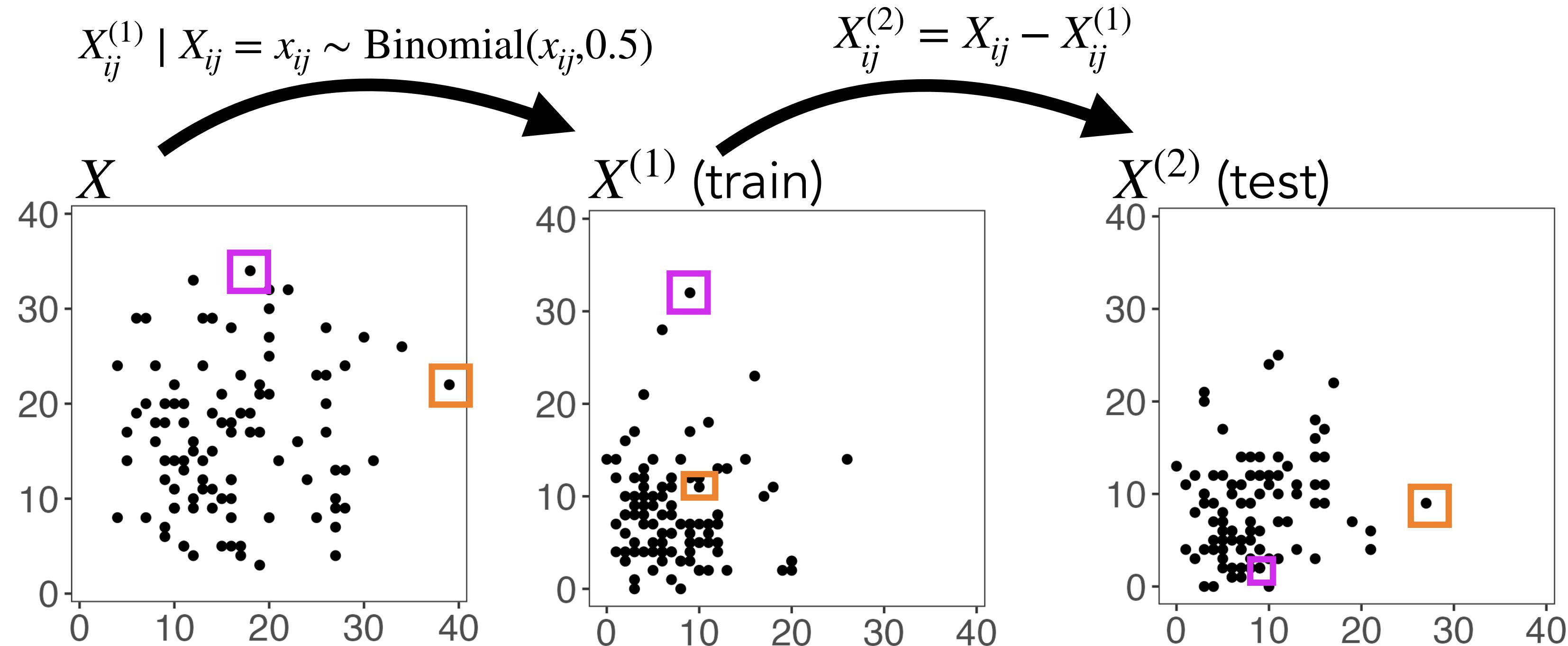
**Step 1:** thin observations into train/test.

# Thinning avoids the pitfall of sample splitting on our motivating examples



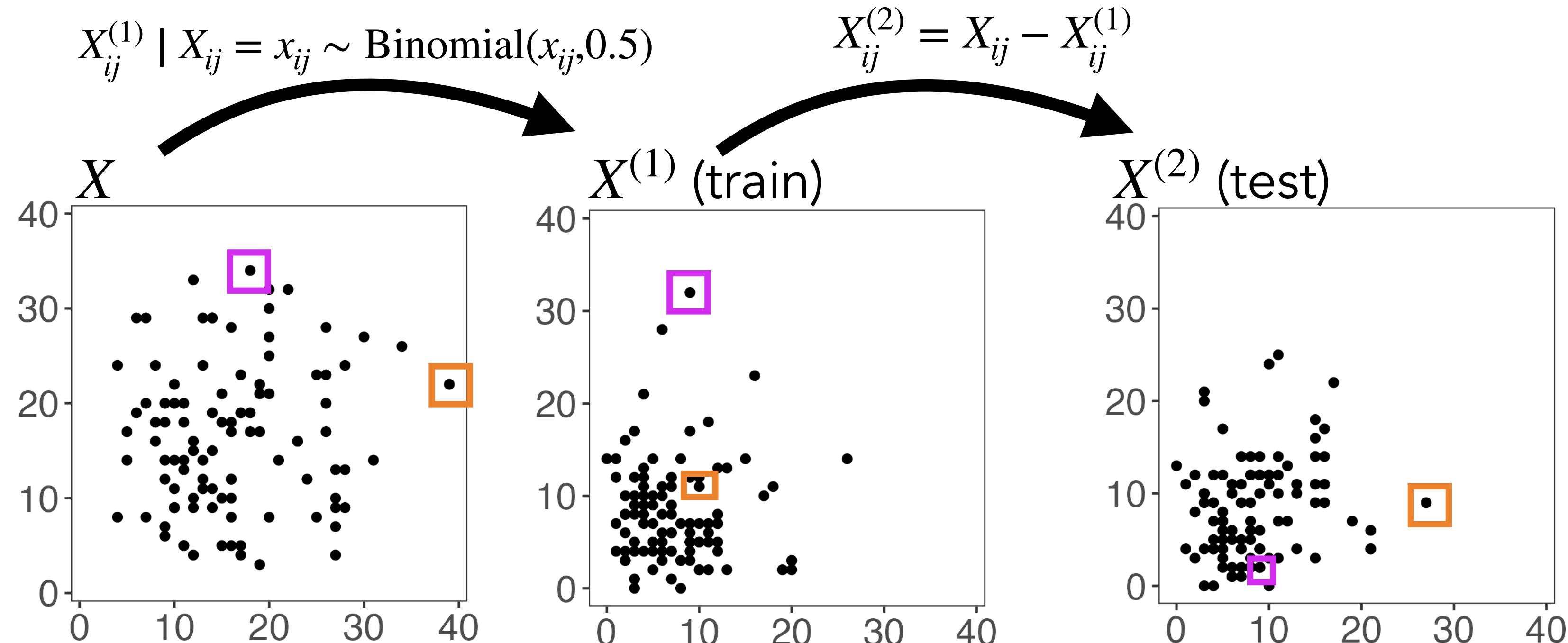
**Step 1:** thin observations into train/test.

# Thinning avoids the pitfall of sample splitting on our motivating examples



**Step 1:** thin observations into train/test.

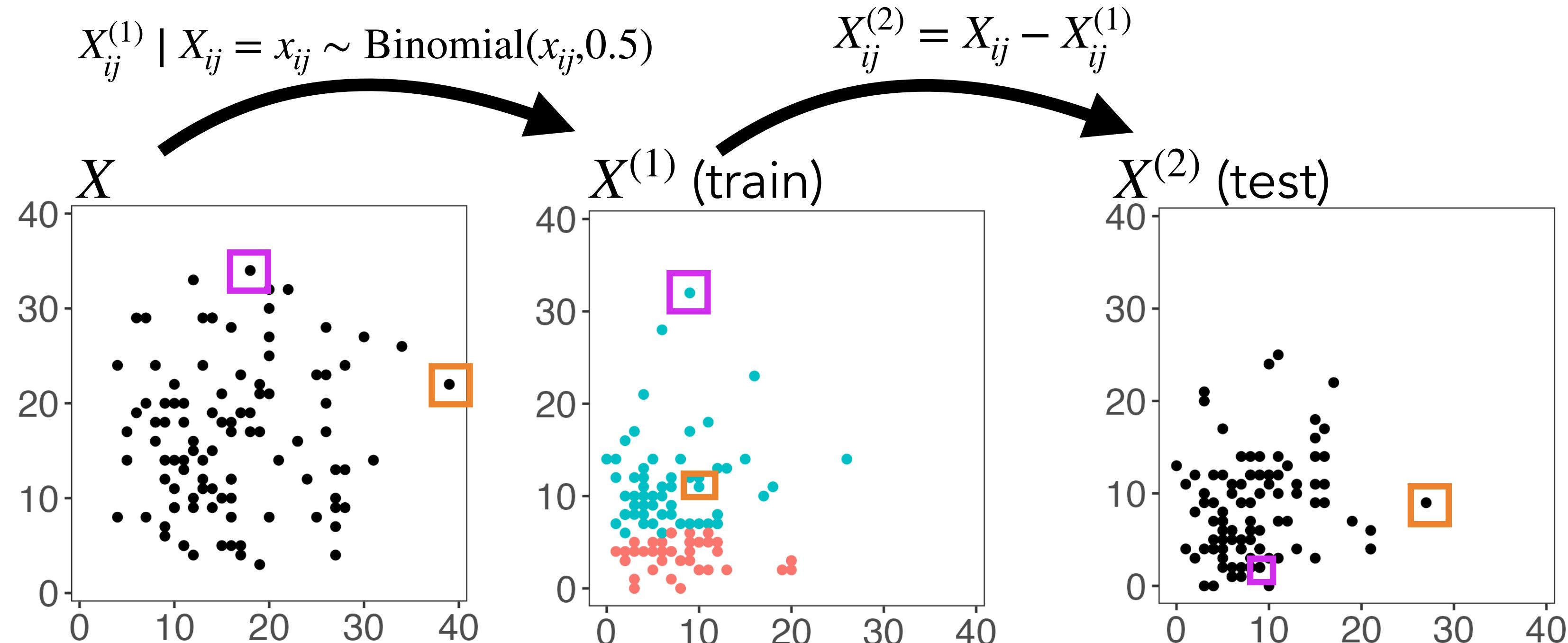
# Thinning avoids the pitfall of sample splitting on our motivating examples



**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

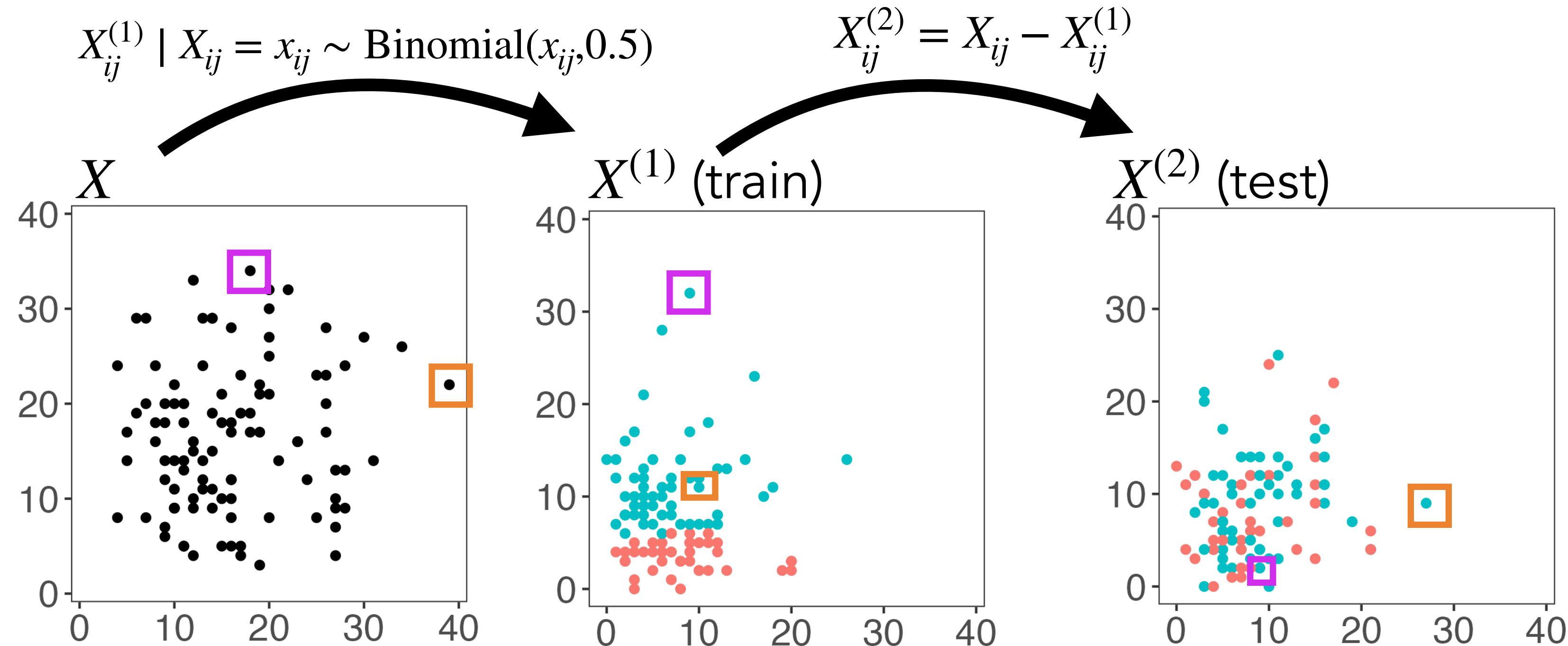
# Thinning avoids the pitfall of sample splitting on our motivating examples



**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

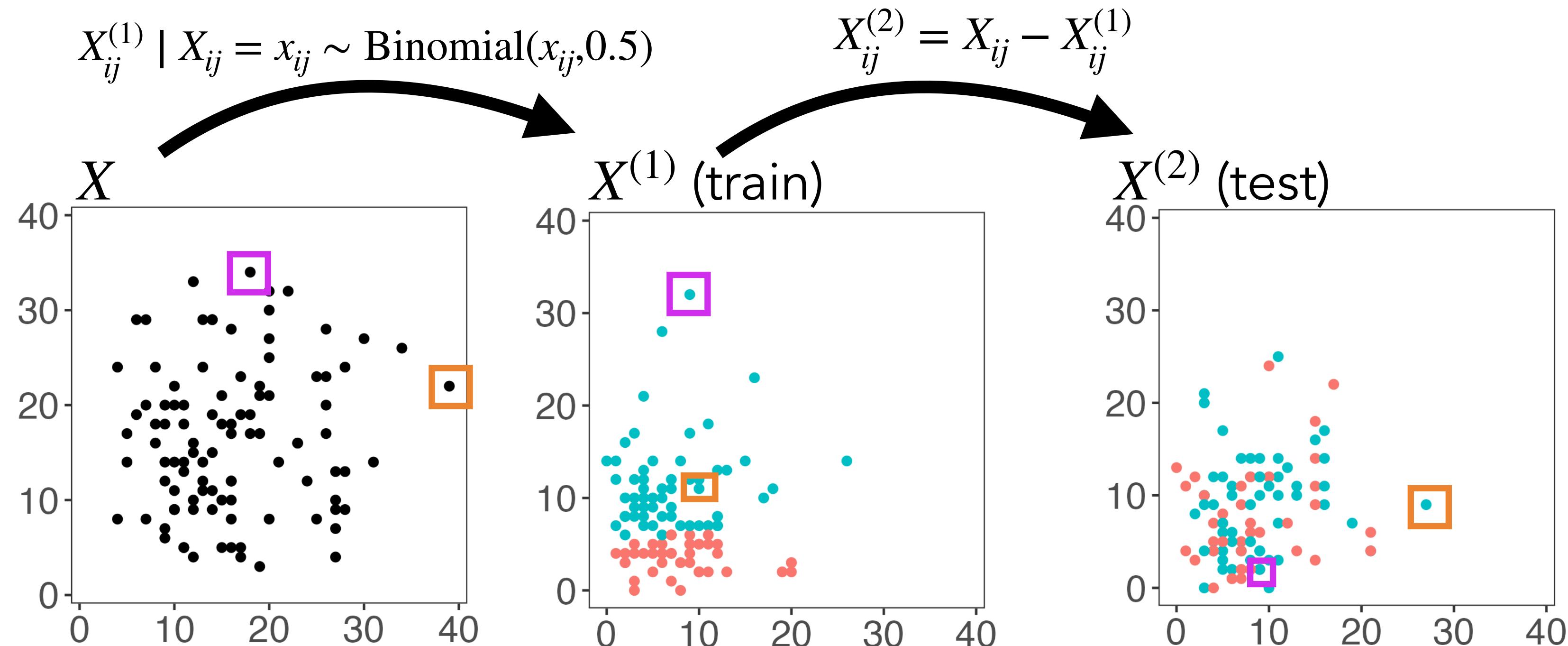
# Thinning avoids the pitfall of sample splitting on our motivating examples



**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

# Thinning avoids the pitfall of sample splitting on our motivating examples

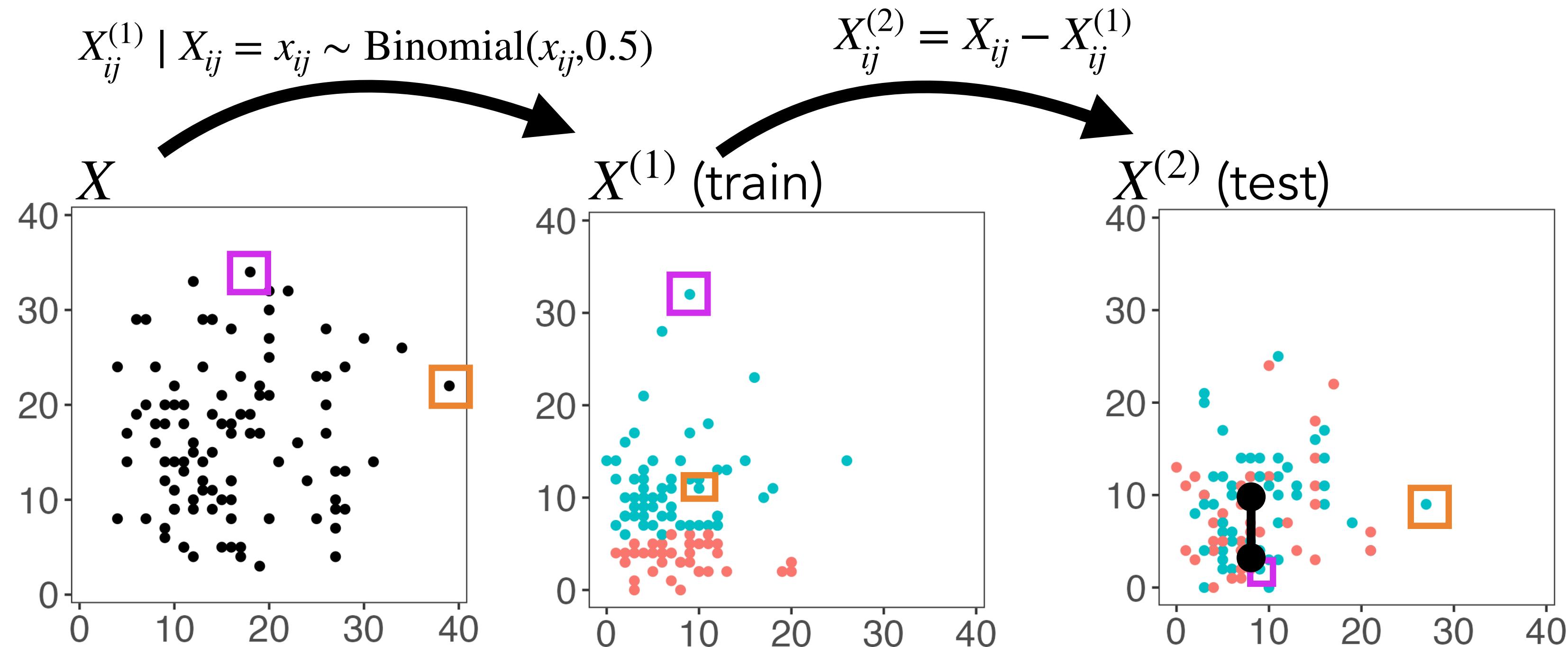


**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

**Step 3:** test for difference in means on test set.

# Thinning avoids the pitfall of sample splitting on our motivating examples

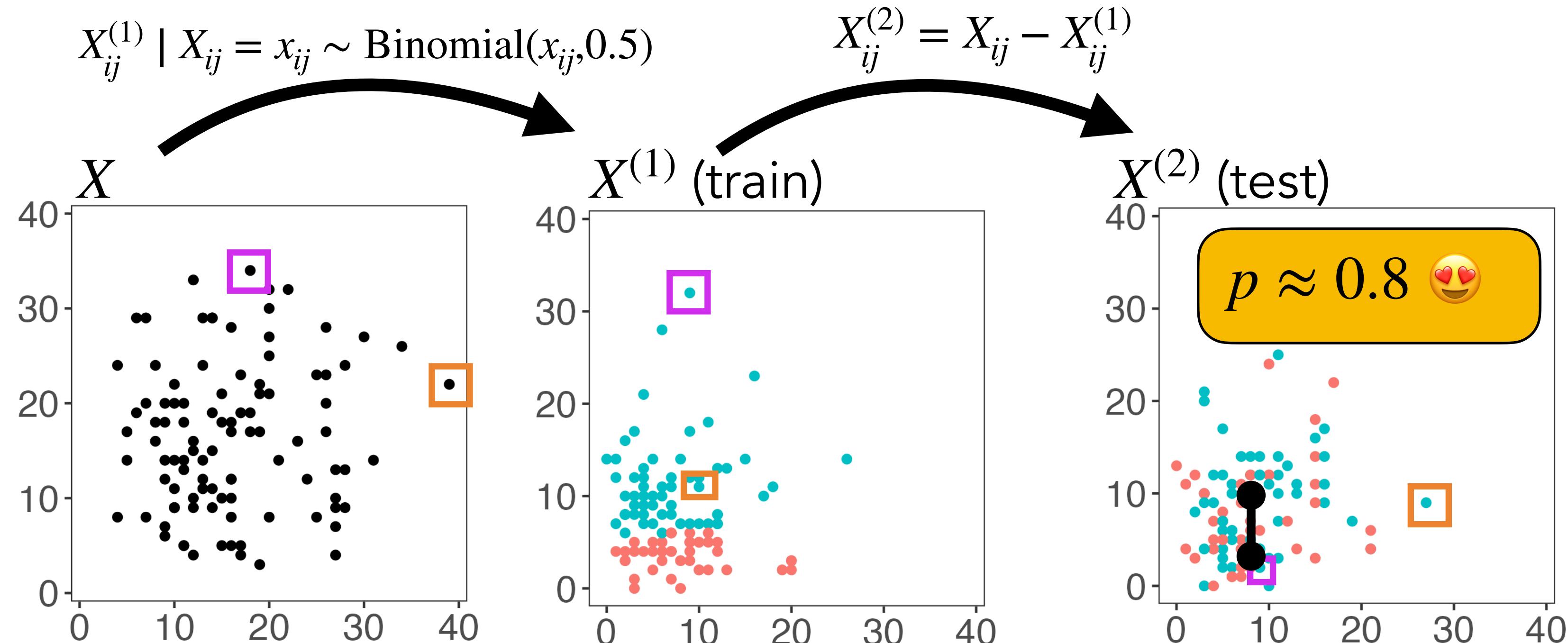


**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

**Step 3:** test for difference in means on test set.

# Thinning avoids the pitfall of sample splitting on our motivating examples

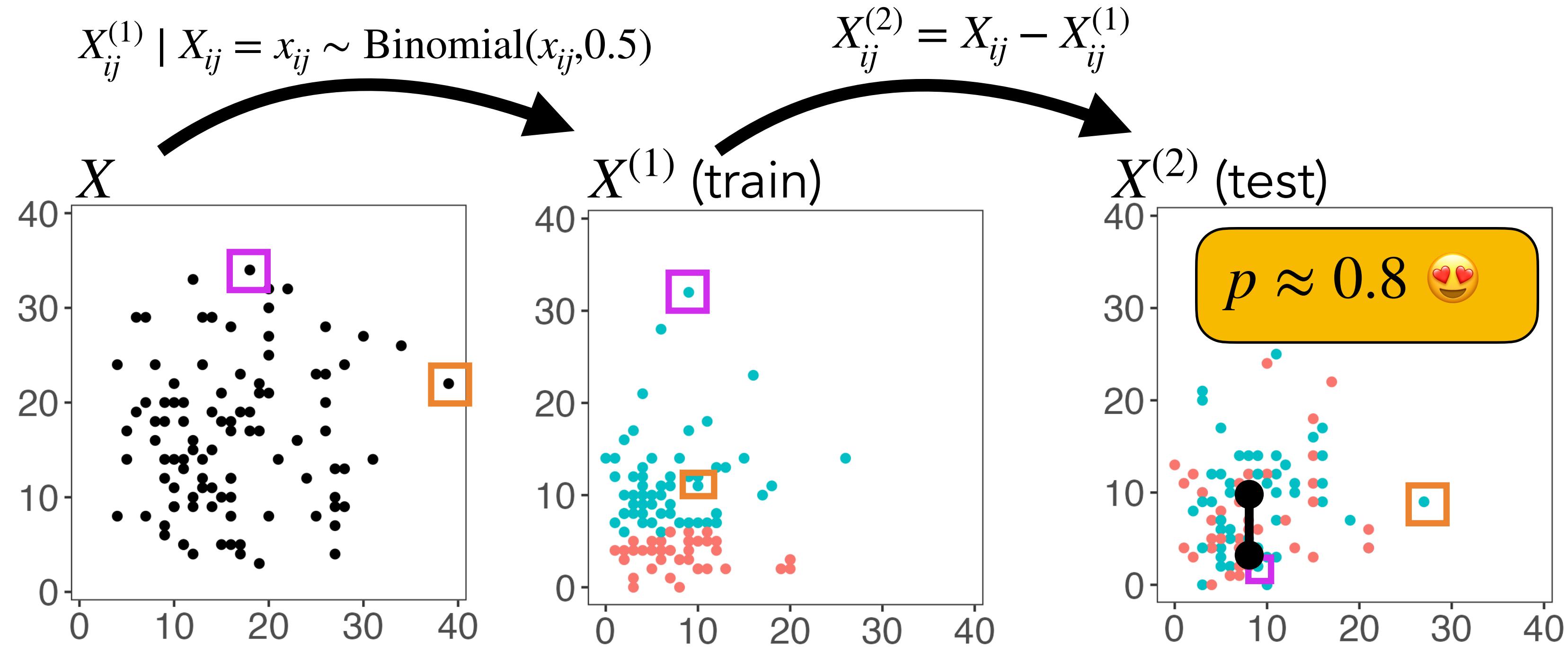


**Step 1:** thin observations into train/test.

**Step 2:** cluster the training set.

**Step 3:** test for difference in means on test set.

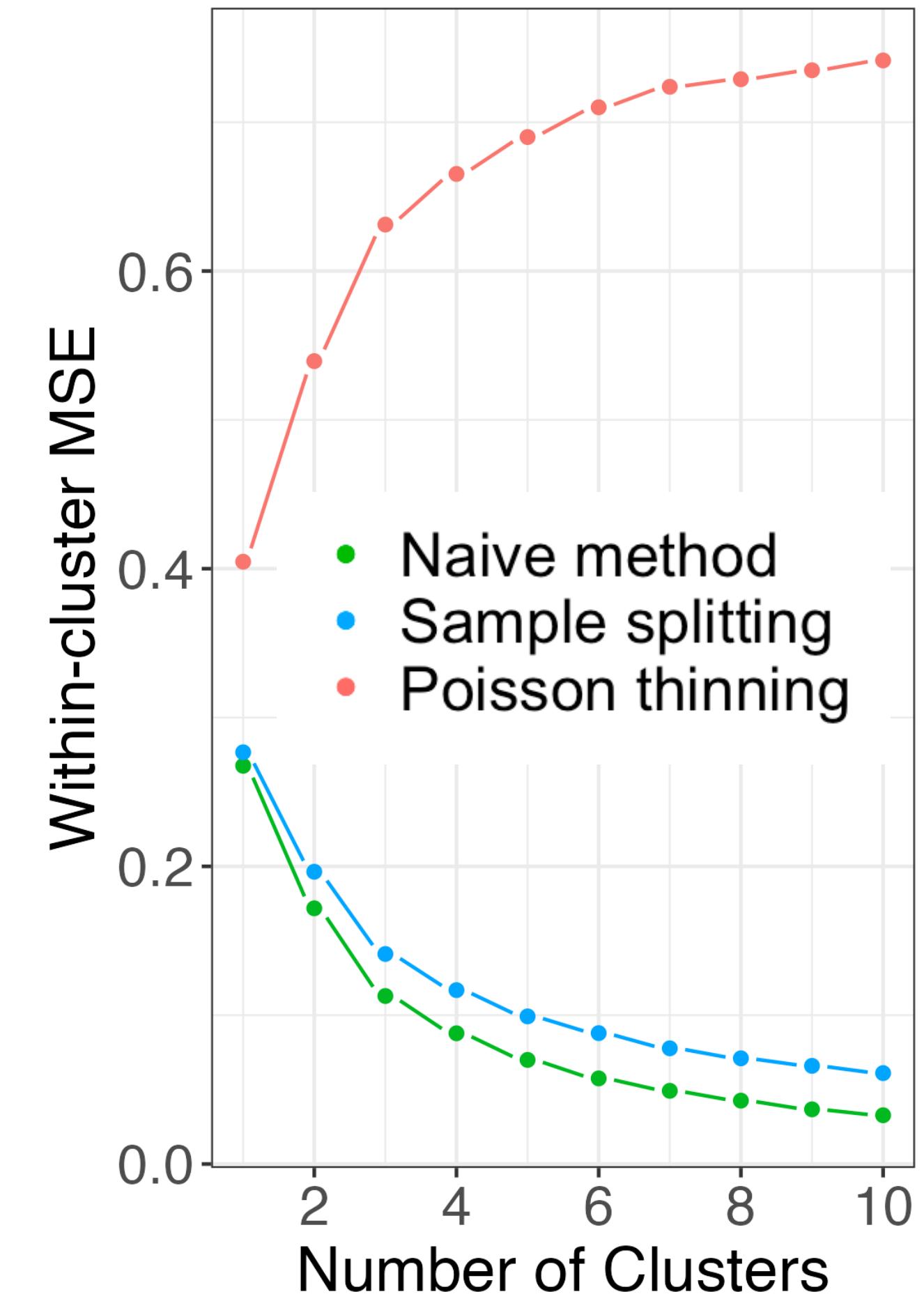
# Thinning avoids the pitfall of sample splitting on our motivating examples



**Step 1:** thin observations into train/test.

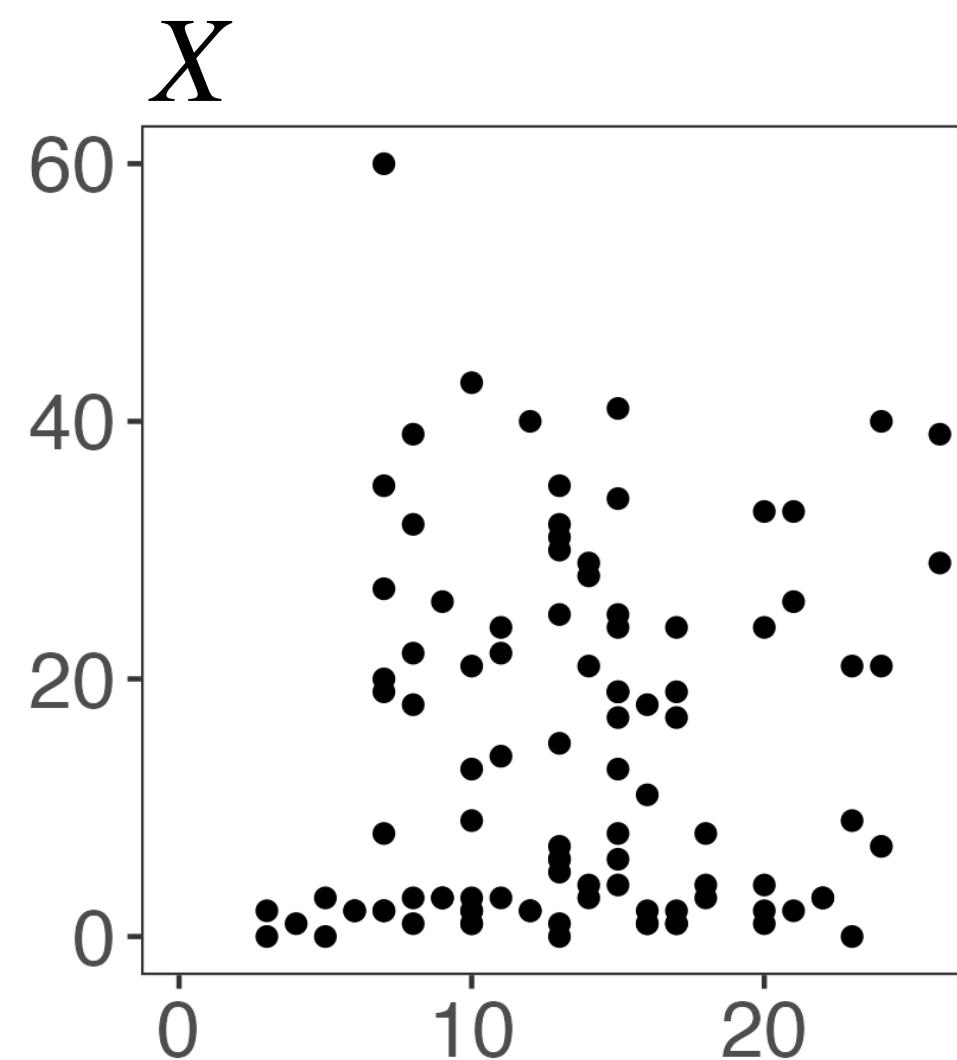
**Step 2:** cluster the training set.

**Step 3:** test for difference in means on test set.



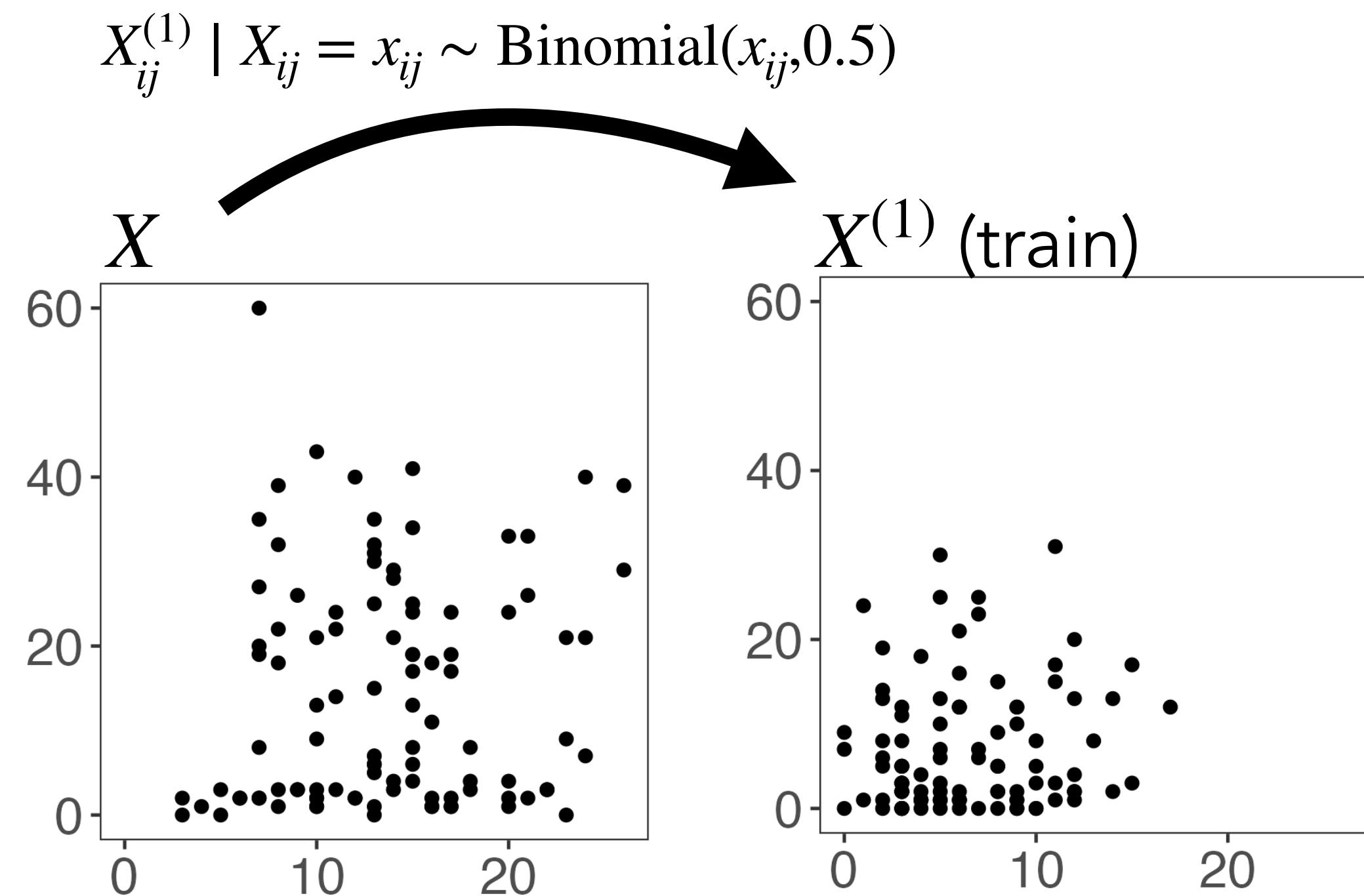
Thinning avoids the pitfall of sample splitting on our motivating examples

---



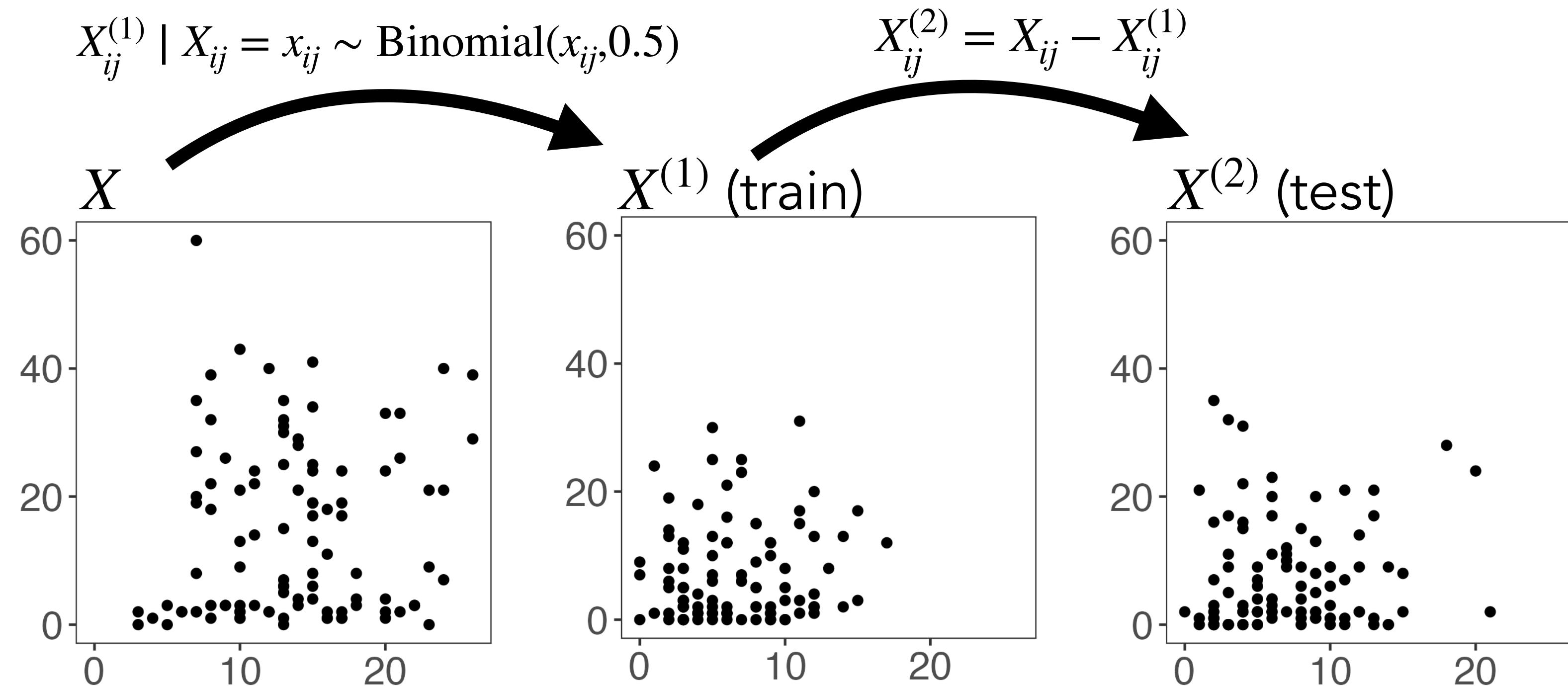
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



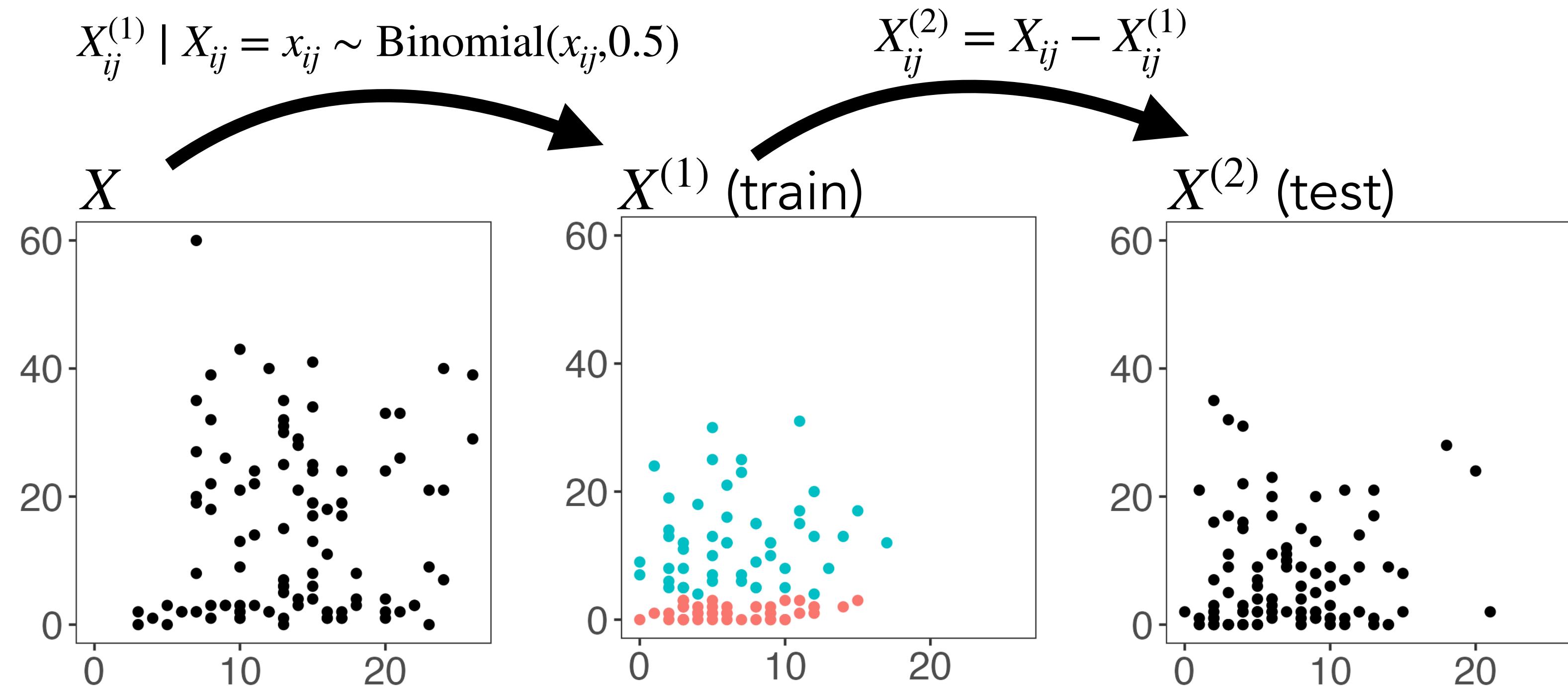
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



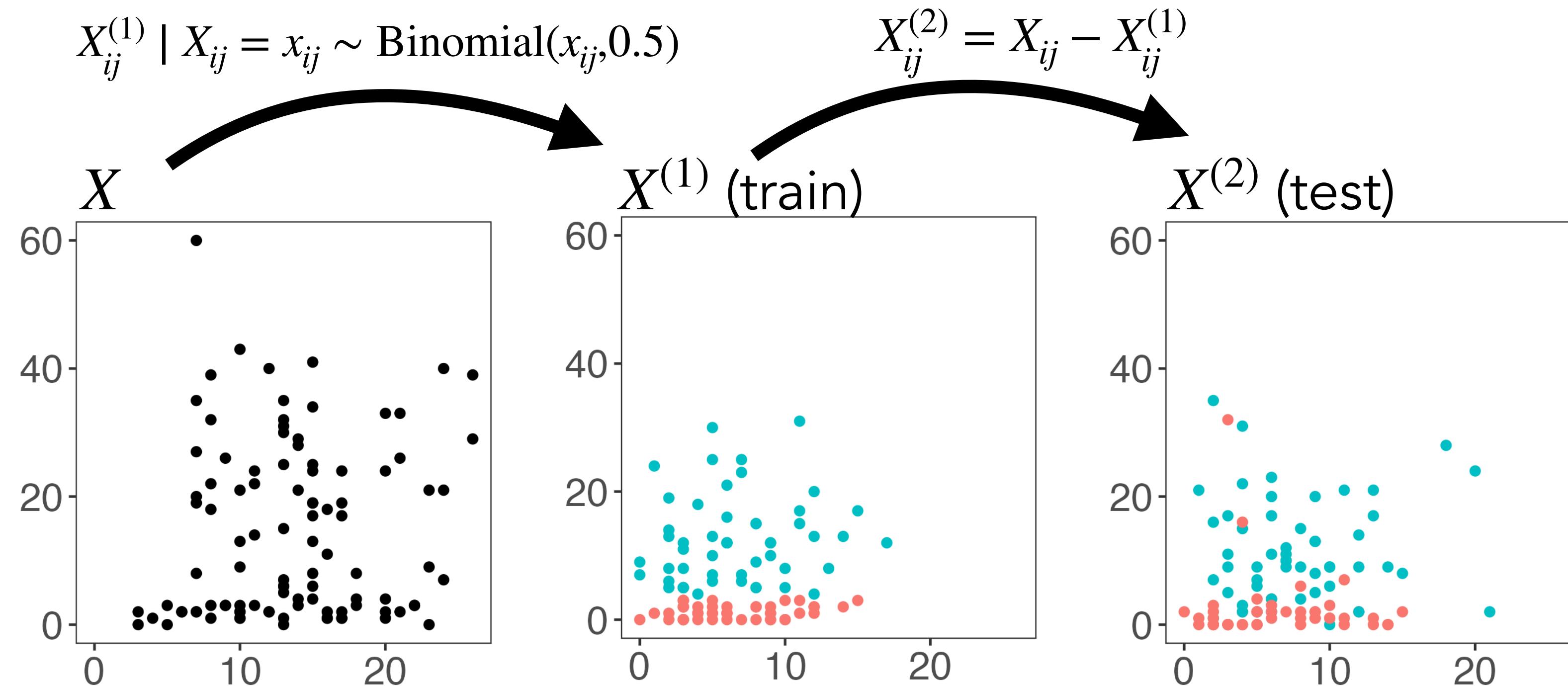
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



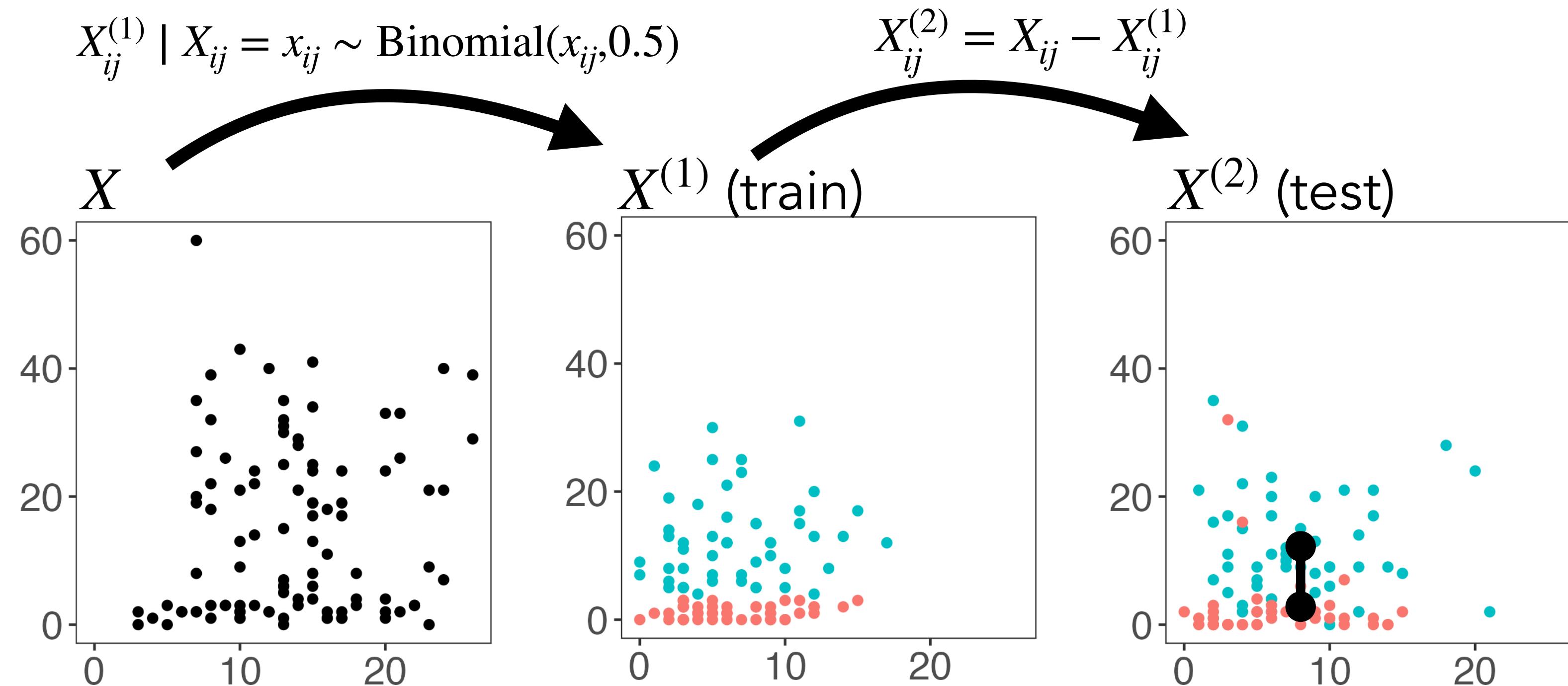
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



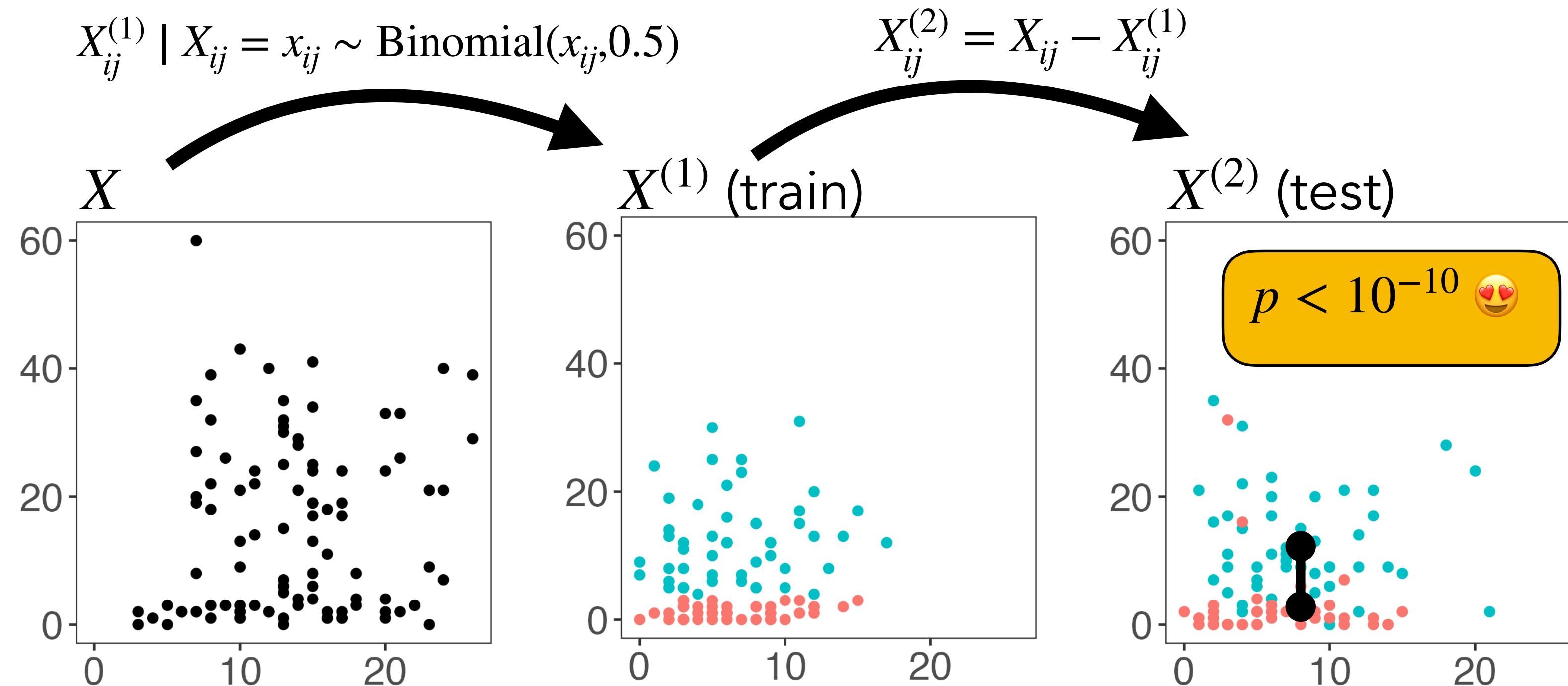
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



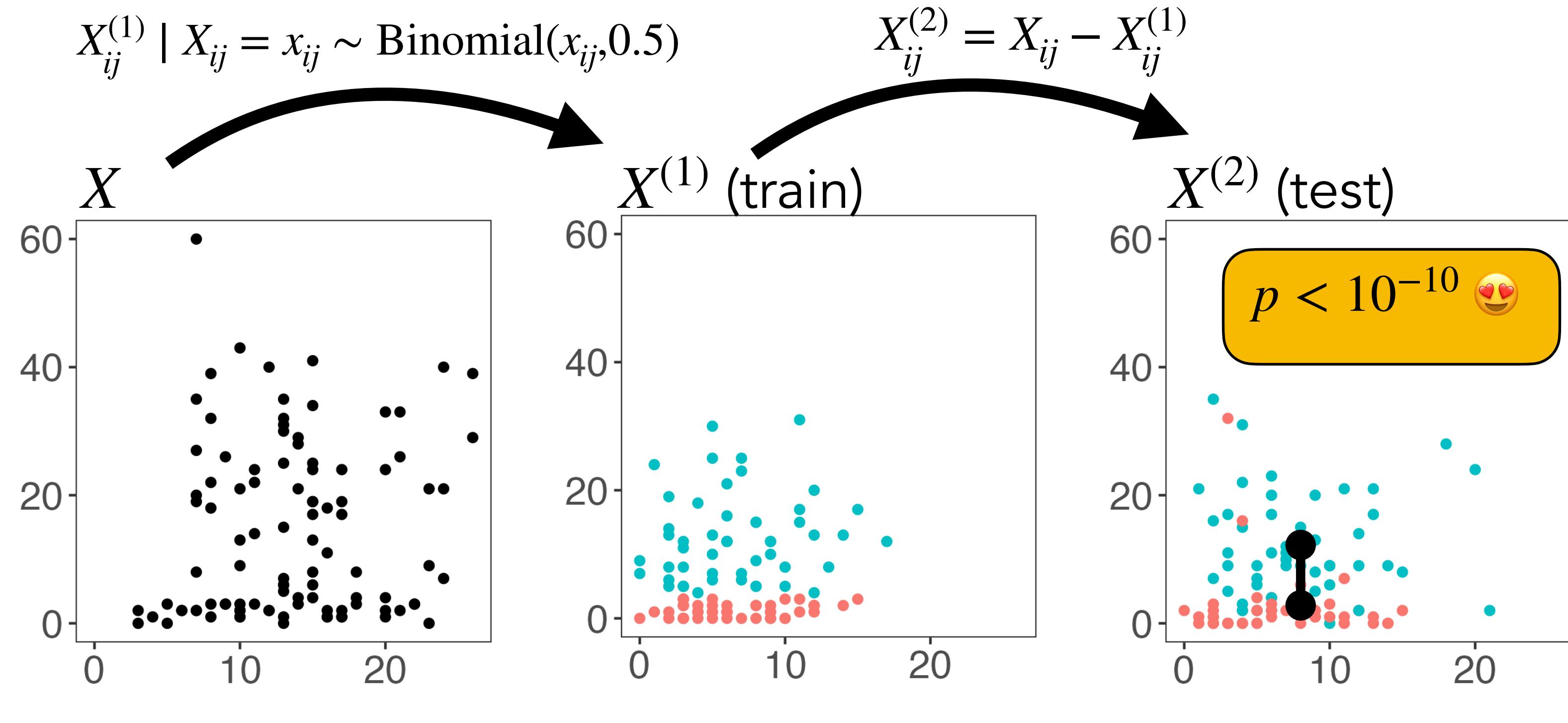
$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples

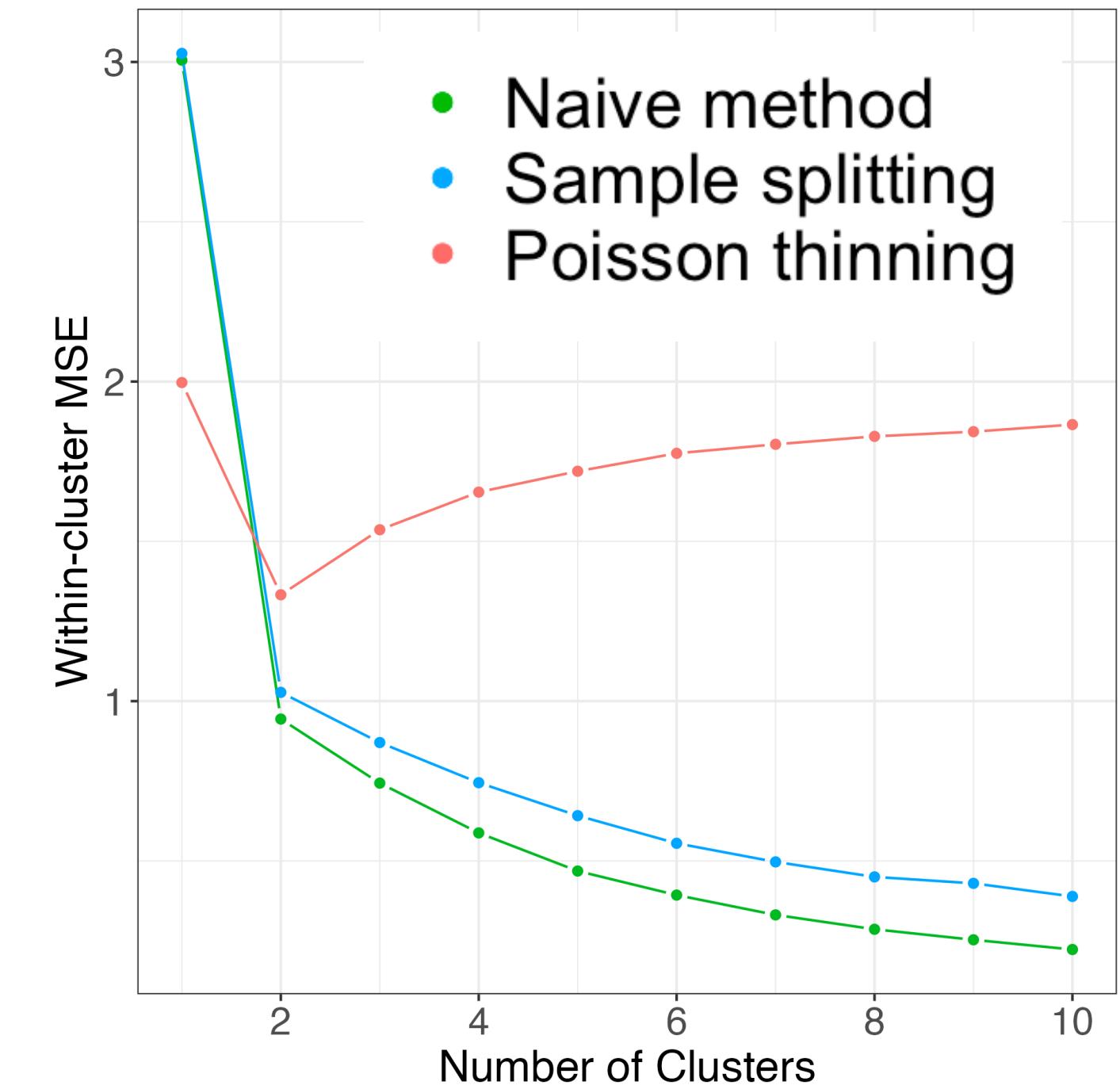


$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$

# Thinning avoids the pitfall of sample splitting on our motivating examples



$$X_{i2} \sim \begin{cases} \text{Poisson}(3) & \text{if } i \leq 50 \\ \text{Poisson}(25) & \text{if } i > 50 \end{cases}$$



# Poisson thinning is useful in the analysis of single-cell RNA sequencing data

*Biostatistics* (2022) **00**, 00, pp. 1–18  
<https://doi.org/10.1093/biostatistics/kxac047>



## Inference after latent variable estimation for single-cell RNA sequencing data

ANNA NEUFELD\*

*Department of Statistics, University of Washington, Seattle, WA 98195, USA*  
aneufeld@uw.edu

LUCY L. GAO

*Department of Statistics, University of British Columbia, BC V6T 1Z4, Canada*

JOSHUA POPP

*Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA*

ALEXIS BATTLE

*Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA and Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218, USA*

DANIELA WITTEN

*Department of Statistics, University of Washington, Seattle, WA 98195, USA and Department of Biostatistics, University of Washington, Seattle, WA 98195, USA*

R package and tutorials:

<https://anna-neufeld.github.io/countssplit/>

# But generalizations of Poisson thinning are needed, even if we only want to study scRNA-seq data

Choudhary and Satija *Genome Biology* (2022) 23:27  
<https://doi.org/10.1186/s13059-021-02584-9>

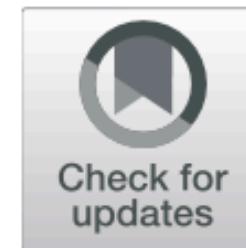
Genome Biology

RESEARCH

Open Access

## Comparison and evaluation of statistical error models for scRNA-seq

Saket Choudhary<sup>1</sup> and Rahul Satija<sup>1,2\*</sup> 



**Results:** Here, we analyze 59 scRNA-seq datasets that span a wide range of technologies, systems, and sequencing depths in order to evaluate the performance of different error models. We find that while a Poisson error model appears appropriate for sparse datasets, we observe clear evidence of overdispersion for genes with sufficient sequencing depth in all biological systems, necessitating the use of a negative binomial model. Moreover, we find that the degree of overdispersion varies widely across datasets, systems, and gene abundances, and argues for a data-driven approach for parameter estimation.

Thinning approaches have the potential to be useful in a wide variety of settings, beyond our motivating example

---

Thinning approaches have the potential to be useful in a wide variety of settings, beyond our motivating example

---

1. Any unsupervised setting, where sample splitting is not an option.

# Thinning approaches have the potential to be useful in a wide variety of settings, beyond our motivating example

---

1. Any unsupervised setting, where sample splitting is not an option.
2. Supervised settings where sample splitting is unsatisfying:
  - Fixed-X regression settings.
  - Non-IID data.
  - Data with outliers or influential points.

# Thinning approaches have the potential to be useful in a wide variety of settings, beyond our motivating example

---

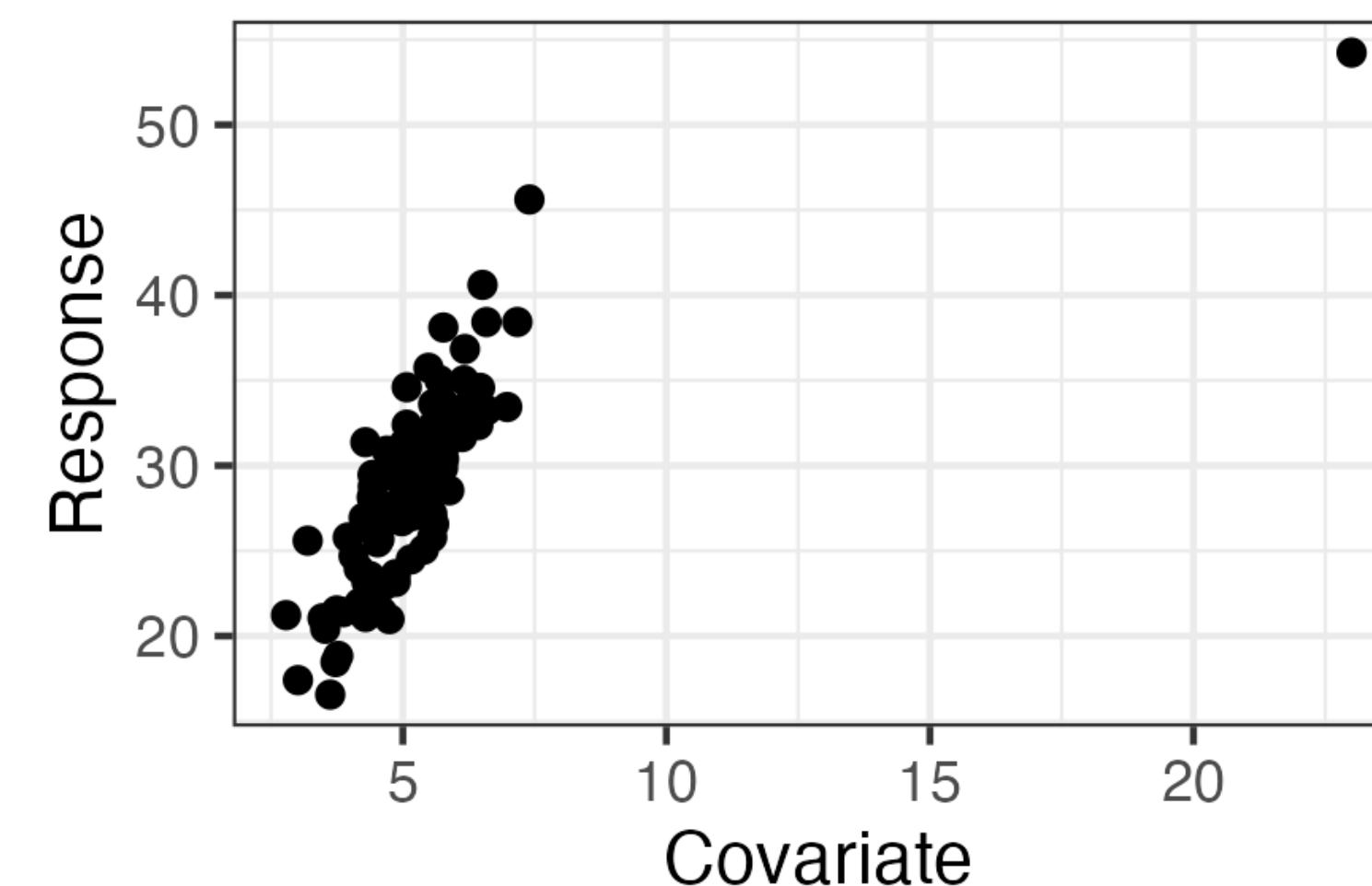
1. Any unsupervised setting, where sample splitting is not an option.
  
2. Supervised settings where sample splitting is unsatisfying:
  - Fixed-X regression settings.
  
  - Non-IID data.
  
  - Data with outliers or influential points.



# Thinning approaches have the potential to be useful in a wide variety of settings, beyond our motivating example

---

1. Any unsupervised setting, where sample splitting is not an option.
  
2. Supervised settings where sample splitting is unsatisfying:
  - Fixed-X regression settings.
  - Non-IID data.
  - Data with outliers or influential points.



# Outline

---

1. Motivation: settings where sample splitting doesn't work
2. Poisson thinning
3. **Data thinning**
4. Application to single-cell RNA sequencing data
5. Ongoing work

## What did we like about Poisson thinning?

---

We split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

## What did we like about Poisson thinning?

---

We split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

Can we achieve these same properties when  $X$  is not Poisson?

## Convolution-closed distributions

---

A family of distributions  $F_\lambda$  is “convolution-closed” in parameter  $\lambda$  if

- $X' \sim F_{\lambda_1}$
- $X'' \sim F_{\lambda_2}$
- $X' \perp\!\!\!\perp X''$

together imply that

$$X' + X'' \sim F_{\lambda_1 + \lambda_2}.$$

# Convolution-closed distributions

A family of distributions  $F_\lambda$  is “convolution-closed” in parameter  $\lambda$  if

- $X' \sim F_{\lambda_1}$
- $X'' \sim F_{\lambda_2}$
- $X' \perp\!\!\!\perp X''$

together imply that

$$X' + X'' \sim F_{\lambda_1 + \lambda_2}.$$

Distribution	Convolution-closed in:
$X \sim \text{Poisson}(\lambda)$	$\lambda$
$X \sim N(\mu, \sigma^2)$	$(\mu, \sigma^2)$
$X \sim \text{NegativeBinomial}(\mu, b)$	$(\mu, b)$
$X \sim \text{Gamma}(\alpha, \beta)$	$\alpha$ , if $\beta$ is fixed
$X \sim \text{Binomial}(r, p)$	$r$ , if $p$ is fixed
$X \sim N_k(\mu, \Sigma)$ .	$(\mu, \Sigma)$ .
$X \sim \text{Multinomial}_k(r, p)$	$r$ , if $p$ is fixed
$X \sim \text{Wishart}_p(n, \Sigma)$	$n$ , if $p$ and $\Sigma$ are fixed.

# Data thinning for convolution-closed distributions

---

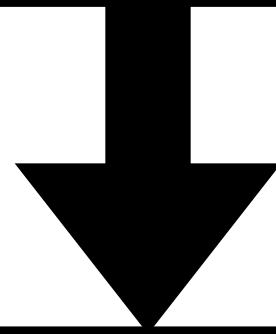
# Data thinning for convolution-closed distributions

---

We observe realization  $x$  from  $X \sim F_\lambda$ .

## Data thinning for convolution-closed distributions

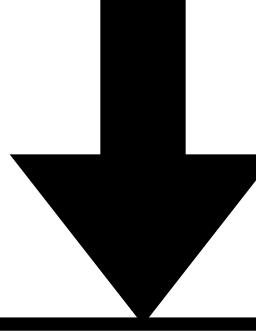
We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim F_\lambda$ .

## Data thinning for convolution-closed distributions

We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .

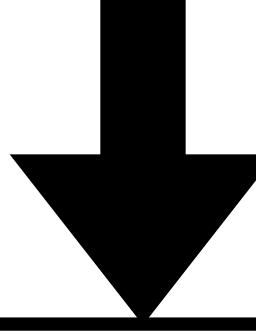


If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

We observe realization  $x$  from  $X \sim F_\lambda$ .

# Data thinning for convolution-closed distributions

We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .



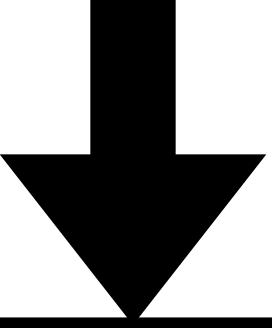
We observe realization  $x$  from  $X \sim F_\lambda$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

# Data thinning for convolution-closed distributions

We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim F_\lambda$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

Let  $G_{\epsilon,x}$  be the conditional distribution of  $X' | X = x$ .

# Data thinning for convolution-closed distributions

We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

We observe realization  $x$  from  $X \sim F_\lambda$ .

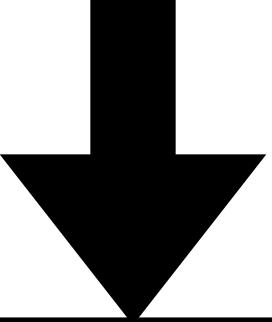
Can we work backwards to recover  $x'$  and  $x''$ ?

Draw  $X^{(1)}$  from  $G_{\epsilon,x}$ . Let  $X^{(2)} := X - X^{(1)}$ .

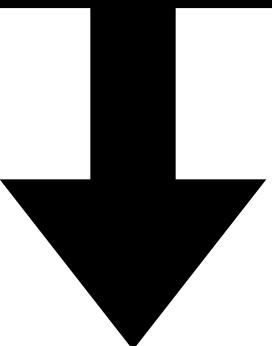
Let  $G_{\epsilon,x}$  be the conditional distribution of  $X' | X = x$ .

# Data thinning for convolution-closed distributions

We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim F_{\epsilon\lambda}$ ,  $X'' \sim F_{(1-\epsilon)\lambda}$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim F_\lambda$ .



Draw  $X^{(1)}$  from  $G_{\epsilon,x}$ . Let  $X^{(2)} := X - X^{(1)}$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

Let  $G_{\epsilon,x}$  be the conditional distribution of  $X' | X = x$ .

**Theorem:**

$X^{(1)} \sim F_{\epsilon\lambda}$ ,  $X^{(2)} \sim F_{(1-\epsilon)\lambda}$ ,  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

# Data thinning for the Poisson distribution

---

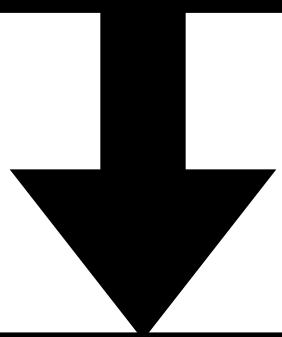
## Data thinning for the Poisson distribution

---

We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .

## Data thinning for the Poisson distribution

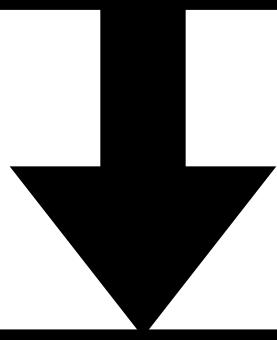
We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .

## Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .

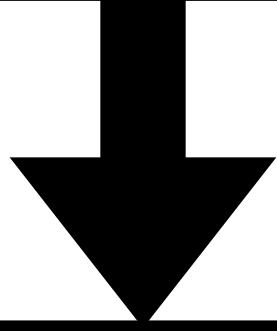


If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .

# Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



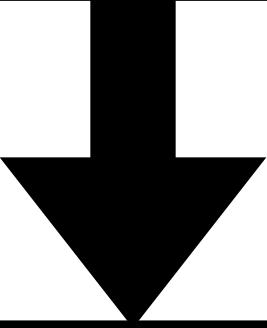
We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

# Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .

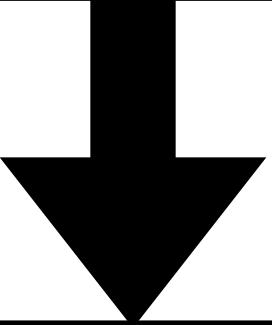
If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

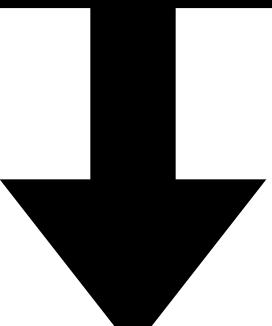
The conditional distribution of  $X' | X = x$  is Binomial( $x, \epsilon$ ).

# Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .



Draw  $X^{(1)}$  from  $\text{Binomial}(x, \epsilon)$ . Let  $X^{(2)} := X - X^{(1)}$ .

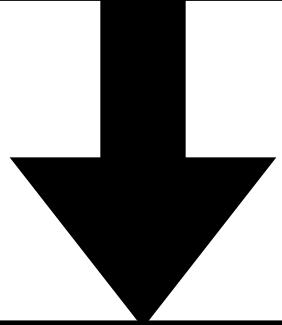
If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

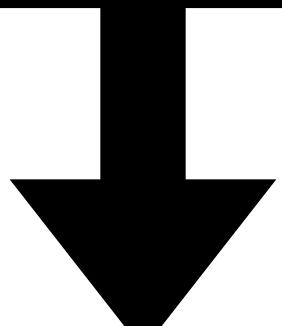
The conditional distribution of  $X' | X = x$  is  $\text{Binomial}(x, \epsilon)$ .

# Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .



Draw  $X^{(1)}$  from  $\text{Binomial}(x, \epsilon)$ . Let  $X^{(2)} := X - X^{(1)}$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

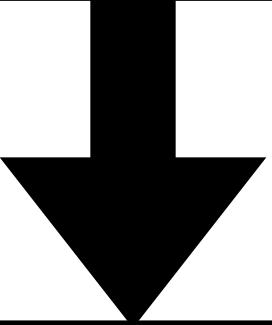
The conditional distribution of  $X' | X = x$  is  $\text{Binomial}(x, \epsilon)$ .

**Theorem:**

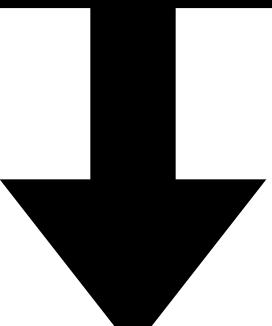
$X^{(1)} \sim \text{Pois}(\epsilon\lambda)$ ,  $X^{(2)} \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

# Data thinning for the Poisson distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{Pois}(\epsilon\lambda)$ ,  $X'' \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{Poisson}(\lambda)$ .



Draw  $X^{(1)}$  from  $\text{Binomial}(x, \epsilon)$ . Let  $X^{(2)} := X - X^{(1)}$ .

## Theorem:

$X^{(1)} \sim \text{Pois}(\epsilon\lambda)$ ,  $X^{(2)} \sim \text{Pois}((1 - \epsilon)\lambda)$ ,  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

The conditional distribution of  $X' | X = x$  is  $\text{Binomial}(x, \epsilon)$ .

We have recovered Poisson thinning!

# Data thinning for the negative binomial distribution

---

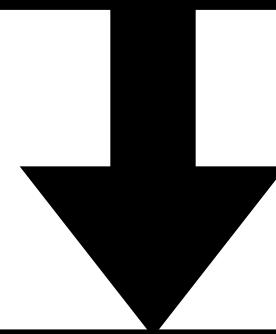
# Data thinning for the negative binomial distribution

---

We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .

## Data thinning for the negative binomial distribution

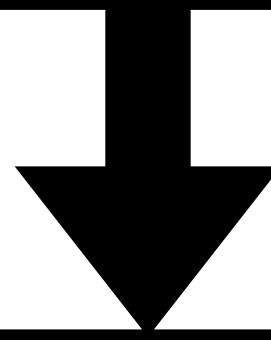
We know  $x$  could have arisen as  $x' + x''$ , where  
 $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .

## Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .

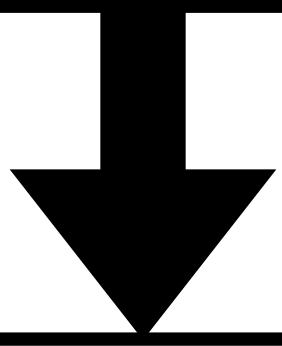


If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .

# Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



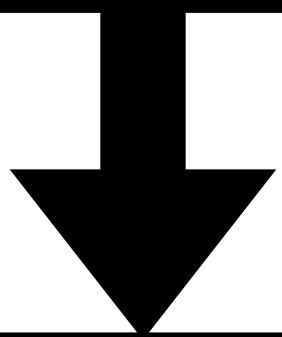
We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

# Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .

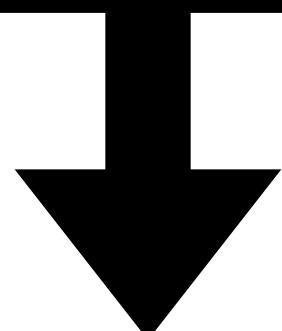
If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

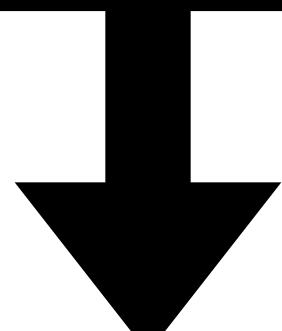
The conditional distribution of  $X' | X = x$  is BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).

# Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .



Draw  $X^{(1)}$  from BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).  
Let  $X^{(2)} := X - X^{(1)}$ .

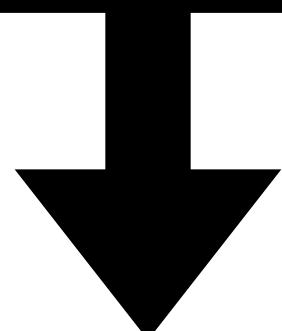
If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

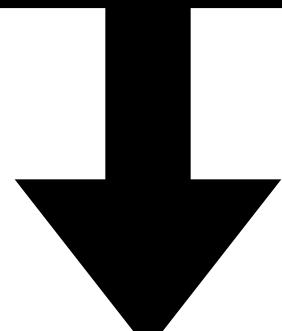
The conditional distribution of  $X' | X = x$  is BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).

# Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .



Draw  $X^{(1)}$  from BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).  
Let  $X^{(2)} := X - X^{(1)}$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

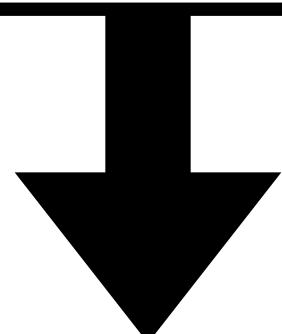
The conditional distribution of  $X' | X = x$  is BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).

## Theorem:

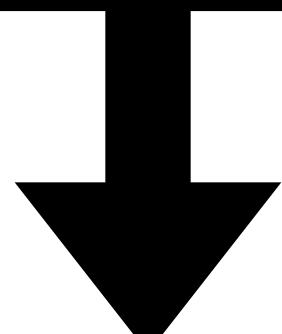
$X^{(1)} \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X^{(2)} \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

# Data thinning for the negative binomial distribution

We know  $x$  could have arisen as  $x' + x''$ , where  $X' \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X'' \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim \text{NB}(\mu, b)$ .



Draw  $X^{(1)}$  from BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).  
Let  $X^{(2)} := X - X^{(1)}$ .

If we had observed  $x'$  and  $x''$ , we would have satisfied our goal of data thinning!

Can we work backwards to recover  $x'$  and  $x''$ ?

The conditional distribution of  $X' | X = x$  is BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).

## Theorem:

$X^{(1)} \sim \text{NB}(\epsilon\mu, \epsilon b)$ ,  $X^{(2)} \sim \text{NB}((1 - \epsilon)\mu, (1 - \epsilon)b)$ ,  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

*This is a new result!*

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

---

Distribution of  $X$ :

Draw  $X^{(1)} \mid X = x$  from  
 $G_{\epsilon,x}$ , where  $G_{\epsilon,x}$  is:

Poisson( $\lambda$ )

Binomial( $x, \epsilon$ )

Distribution of  $X^{(1)}$ :

Poisson( $\epsilon\lambda$ )

Distribution of  $X^{(2)}$ ,

where  $X^{(2)} = X - X^{(1)}$ :

Poisson( $(1 - \epsilon)\lambda$ )

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

---

Distribution of $X$ :	Draw $X^{(1)}   X = x$ from $G_{\epsilon,x}$ , where $G_{\epsilon,x}$ is:	Distribution of $X^{(1)}$ :	Distribution of $X^{(2)}$ , where $X^{(2)} = X - X^{(1)}$ :
Poisson( $\lambda$ )	Binomial( $x, \epsilon$ )	Poisson( $\epsilon\lambda$ )	Poisson( $(1 - \epsilon)\lambda$ )

## Related work on Poisson thinning:

- Sarkar and Stephens, 2021, Nature Genetics.
- Chen et al., 2021, arXiv:2108.03336
- Leiner et al., 2023, JASA
- Neufeld et al., 2022, Biostatistics.
- Oliveira, Lei, and Tibshirani, 2022, arXiv:2212.01943.

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

---

Distribution of $X$ :	Draw $X^{(1)}   X = x$ from $G_{\epsilon,x}$ , where $G_{\epsilon,x}$ is:	Distribution of $X^{(1)}$ :	Distribution of $X^{(2)}$ , where $X^{(2)} = X - X^{(1)}$ :
Poisson( $\lambda$ )	Binomial( $x, \epsilon$ )	Poisson( $\epsilon\lambda$ )	Poisson( $(1 - \epsilon)\lambda$ )
$N(\mu, \sigma^2)$	$N(\epsilon x, \epsilon(1 - \epsilon)\sigma^2)$	$N(\epsilon\mu, \epsilon\sigma^2)$	$N((1 - \epsilon)\mu, (1 - \epsilon)\sigma^2)$

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

Distribution of $X$ :	Draw $X^{(1)}   X = x$ from $G_{\epsilon,x}$ , where $G_{\epsilon,x}$ is:	Distribution of $X^{(1)}$ :	Distribution of $X^{(2)}$ , where $X^{(2)} = X - X^{(1)}$ :
Poisson( $\lambda$ )	Binomial( $x, \epsilon$ )	Poisson( $\epsilon\lambda$ )	Poisson( $(1 - \epsilon)\lambda$ )
$N(\mu, \sigma^2)$	$N(\epsilon x, \epsilon(1 - \epsilon)\sigma^2)$	$N(\epsilon\mu, \epsilon\sigma^2)$	$N((1 - \epsilon)\mu, (1 - \epsilon)\sigma^2)$

## Related work on Gaussian thinning:

- Tian and Taylor, 2018, Annals of Statistics.
- Tian, 2020, Annals of Statistics.
- Rasines and Young, 2022, Biometrika.
- Leiner et al., 2023, JASA
- Oliveira, Lei, and Tibshirani, 2022, arXiv:2111.09447.

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

---

Distribution of $X$ :	Draw $X^{(1)} \mid X = x$ from $G_{\epsilon,x}$ , where $G_{\epsilon,x}$ is:	Distribution of $X^{(1)}$ :	Distribution of $X^{(2)}$ , where $X^{(2)} = X - X^{(1)}$ :
Poisson( $\lambda$ )	Binomial( $x, \epsilon$ )	Poisson( $\epsilon\lambda$ )	Poisson( $(1 - \epsilon)\lambda$ )
$N(\mu, \sigma^2)$	$N(\epsilon x, \epsilon(1 - \epsilon)\sigma^2)$	$N(\epsilon\mu, \epsilon\sigma^2)$	$N((1 - \epsilon)\mu, (1 - \epsilon)\sigma^2)$
NegativeBinomial( $\mu, b$ )	BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).	NegativeBinomial( $\epsilon\mu, \epsilon b$ )	NegativeBinomial( $(1 - \epsilon)\mu, (1 - \epsilon)b$ )

# For many common distributions, the distribution $G_{\epsilon,x}$ has a simple form

---

Distribution of $X$ :	Draw $X^{(1)} \mid X = x$ from $G_{\epsilon,x}$ , where $G_{\epsilon,x}$ is:	Distribution of $X^{(1)}$ :	Distribution of $X^{(2)}$ , where $X^{(2)} = X - X^{(1)}$ :
Poisson( $\lambda$ )	Binomial( $x, \epsilon$ )	Poisson( $\epsilon\lambda$ )	Poisson( $(1 - \epsilon)\lambda$ )
$N(\mu, \sigma^2)$	$N(\epsilon x, \epsilon(1 - \epsilon)\sigma^2)$	$N(\epsilon\mu, \epsilon\sigma^2)$	$N((1 - \epsilon)\mu, (1 - \epsilon)\sigma^2)$
NegativeBinomial( $\mu, b$ )	BetaBinomial( $x, \epsilon b, (1 - \epsilon)b$ ).	NegativeBinomial( $\epsilon\mu, \epsilon b$ )	NegativeBinomial( $(1 - \epsilon)\mu, (1 - \epsilon)b$ )
Binomial( $r, p$ )	Hypergeometric( $\epsilon r, (1 - \epsilon)r, x$ ).	Binomial( $\epsilon r, p$ )	Binomial( $(1 - \epsilon)r, p$ )
Gamma( $\alpha, \beta$ )	$x \cdot \text{Beta}(\epsilon\alpha, (1 - \epsilon)\alpha)$ .	Gamma( $\epsilon\alpha, \beta$ )	Gamma( $(1 - \epsilon)\alpha, \beta$ )
Exponential( $\lambda$ )	$x \cdot \text{Beta}(\epsilon, (1 - \epsilon))$ .	Gamma( $\epsilon, \lambda$ )	Gamma( $(1 - \epsilon), \lambda$ )
$N_k(\mu, \Sigma)$	$N(\epsilon x, \epsilon(1 - \epsilon)\Sigma)$ .	$N_k(\epsilon\mu, \epsilon\Sigma)$	$N_k((1 - \epsilon)\mu, (1 - \epsilon)\Sigma)$
Multinomial $_k(r, p)$	MultivarHypergeom( $x_1, \dots, x_K, \epsilon r$ )	Multinom $_k(\epsilon r, p)$	Multinomial $_k((1 - \epsilon)r, p)$
Wishart $_p(n, \Sigma)$ .	$x^{1/2} Z x^{1/2}$ , where .  $Z \sim \text{MatrixBeta}_p(\epsilon n/2, (1 - \epsilon)n/2)$	Wishart $_p(\epsilon n, \Sigma)$	Wishart $_p((1 - \epsilon)n, \Sigma)$

---

# What if we get a nuisance parameter wrong?

## Negative binomial thinning algorithm

Suppose  $X \sim \text{NegBin}(\mu, b)$ .

Draw

$X^{(1)} \sim \text{BetaBinomial}(x, \epsilon b, (1 - \epsilon)b)$ ,

$X^{(2)} = X - X^{(1)}$ , then:

- 1)  $X^{(1)} \sim \text{NegBin}(\epsilon\mu, \epsilon b)$ .
- 2)  $X^{(2)} \sim \text{NegBin}((1 - \epsilon)\mu, (1 - \epsilon)b)$
- 3)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

# What if we get a nuisance parameter wrong?

## Negative binomial thinning algorithm

Suppose  $X \sim \text{NegBin}(\mu, b)$ .

Draw

$X^{(1)} \sim \text{BetaBinomial}(x, \epsilon \tilde{b}, (1 - \epsilon) \tilde{b})$ ,

$X^{(2)} = X - X^{(1)}$ , then:

- 1)  $X^{(1)} \sim \text{NegBin}(\epsilon \mu, \epsilon b)$ .
- 2)  $X^{(2)} \sim \text{NegBin}((1 - \epsilon) \mu, (1 - \epsilon) b)$
- 3)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

# What if we get a nuisance parameter wrong?

## Negative binomial thinning algorithm

Suppose  $X \sim \text{NegBin}(\mu, b)$ .

Draw

$X^{(1)} \sim \text{BetaBinomial}(x, \epsilon \tilde{b}, (1 - \epsilon) \tilde{b})$ ,

$X^{(2)} = X - X^{(1)}$ , then:

1)  $\cancel{X^{(1)} \sim \text{NegBin}(c\mu, cb)}$ .

2)  $\cancel{X^{(2)} \sim \text{NegBin}((1 - c)\mu, (1 - c)b)}$

3)  $\cancel{X^{(1)} \perp\!\!\!\perp X^{(2)}}$ .

# What if we get a nuisance parameter wrong?

## Negative binomial thinning algorithm

Suppose  $X \sim \text{NegBin}(\mu, b)$ .

Draw

$X^{(1)} \sim \text{BetaBinomial}(x, \epsilon \tilde{b}, (1 - \epsilon) \tilde{b})$ ,

$X^{(2)} = X - X^{(1)}$ , then:

$$1) E[X^{(1)}] = \epsilon \mu.$$

$$2) E[X^{(2)}] = (1 - \epsilon) \mu$$

$$3) \text{Cov}(X^{(1)}, X^{(2)}) = \epsilon(1 - \epsilon) \frac{\mu^2}{b} \left(1 - \frac{b + 1}{\tilde{b} + 1}\right).$$

# What if we get a nuisance parameter wrong?

## Negative binomial thinning algorithm

Suppose  $X \sim \text{NegBin}(\mu, b)$ .

Draw

$X^{(1)} \sim \text{BetaBinomial}(x, \epsilon \tilde{b}, (1 - \epsilon) \tilde{b})$ ,

$X^{(2)} = X - X^{(1)}$ , then:

1)  $E[X^{(1)}] = \epsilon \mu$ .

2)  $E[X^{(2)}] = (1 - \epsilon) \mu$

3)  $\text{Cov}(X^{(1)}, X^{(2)}) = \epsilon(1 - \epsilon) \frac{\mu^2}{b} \left(1 - \frac{b + 1}{\tilde{b} + 1}\right)$ .

Similar results can be derived for other decompositions.

# Data thinning is a simple alternative to sample splitting that can be used in a variety of settings

The screenshot shows a red header with the arXiv logo and navigation links for 'Search...', 'Help | Advanced...'. Below the header, the category 'Statistics > Methodology' is shown, along with the submission date '[Submitted on 18 Jan 2023]'. The main title 'Data thinning for convolution-closed distributions' is displayed in bold. The authors listed are Anna Neufeld, Ameer Dharamshi, Lucy L. Gao, and Daniela Witten. The abstract text describes data thinning as a new approach for splitting observations into independent parts that sum to the original observation, applicable to convolution-closed distributions like Gaussian, Poisson, and binomial.

We propose data thinning, a new approach for splitting an observation into two or more independent parts that sum to the original observation, and that follow the same distribution as the original observation, up to a (known) scaling of a parameter. This proposal is very general, and can be applied to any observation drawn from a "convolution closed" distribution, a class that includes the Gaussian, Poisson, negative binomial, Gamma, and binomial distributions, among others. It is similar in spirit to -- but distinct from, and more easily applicable than -- a recent proposal known as data fission. Data thinning has a number of applications to model selection, evaluation, and inference. For instance, cross-validation via data thinning provides an attractive alternative to the "usual" approach of cross-validation via sample splitting, especially in unsupervised settings in which the latter is not applicable. In simulations and in an application to single-cell RNA-sequencing data, we show that data thinning can be used to validate the results of unsupervised learning approaches, such as k-means clustering and principal components analysis.

R package and tutorials: <https://anna-neufeld.github.io/datathin/>

## Outline

---

1. Motivation: settings where sample splitting doesn't work
2. Poisson thinning
3. Data thinning
4. **Application to single-cell RNA sequencing data**
5. Ongoing work

# How can we validate the results of clustering?

## RESEARCH ARTICLE

HUMAN GENOMICS

### A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

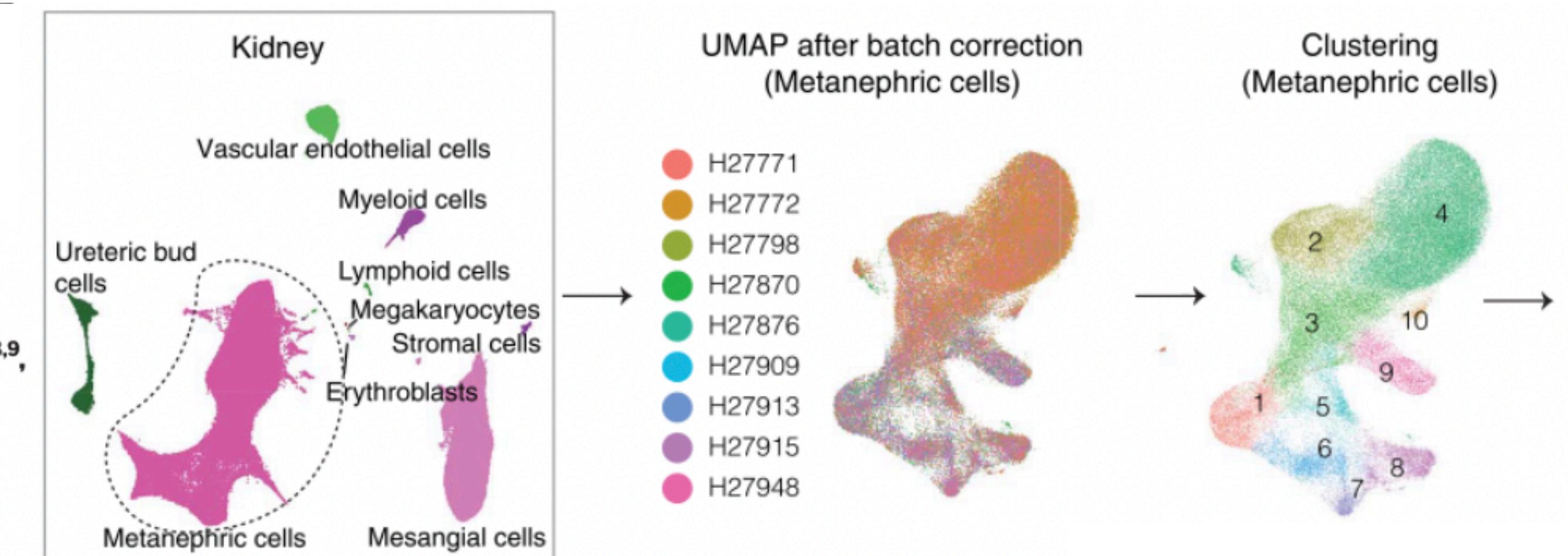
# How can we validate the results of clustering?

## RESEARCH ARTICLE

### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>



# How can we validate the results of clustering?

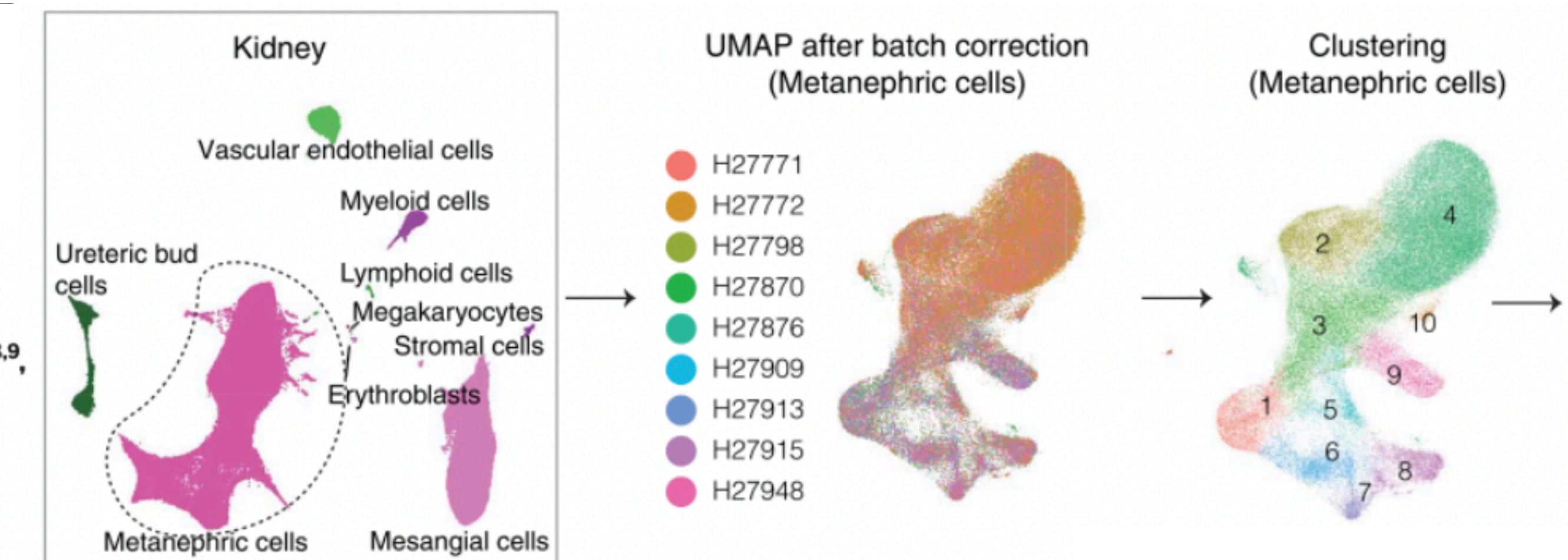
## RESEARCH ARTICLE

### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

“Intradataset cross validation”



# How can we validate the results of clustering?

## RESEARCH ARTICLE

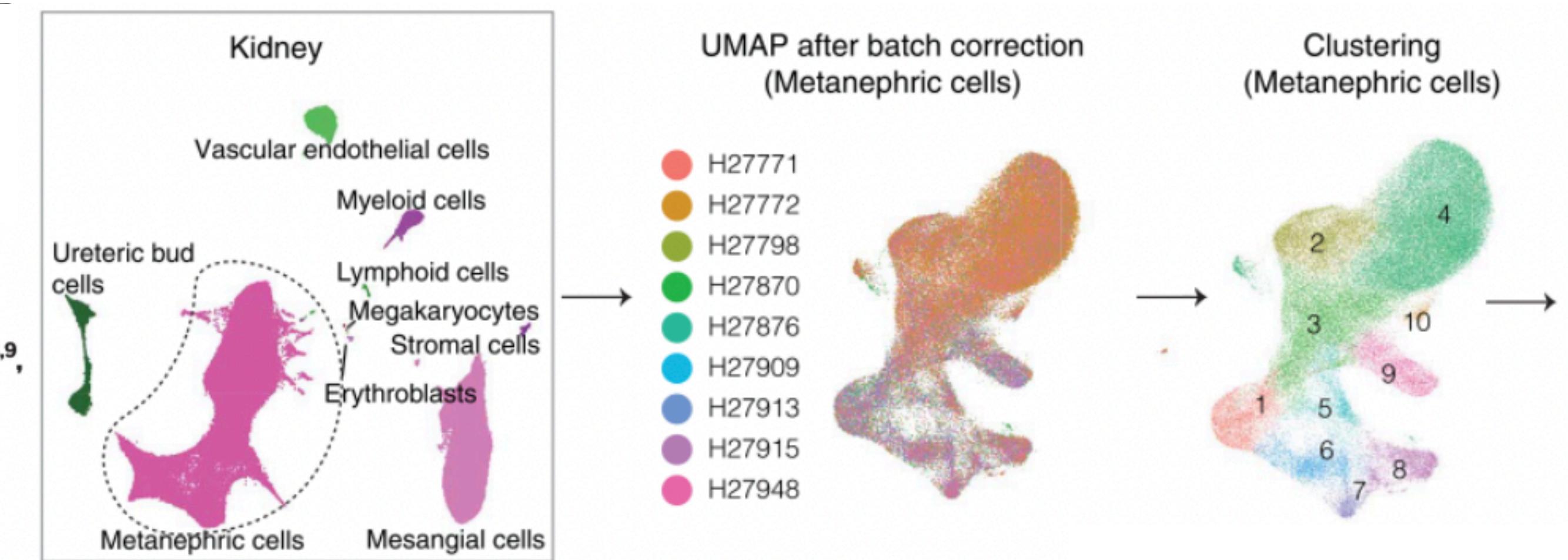
### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

## "Intradataset cross validation"

- Step 1: Cluster cells.



# How can we validate the results of clustering?

## RESEARCH ARTICLE

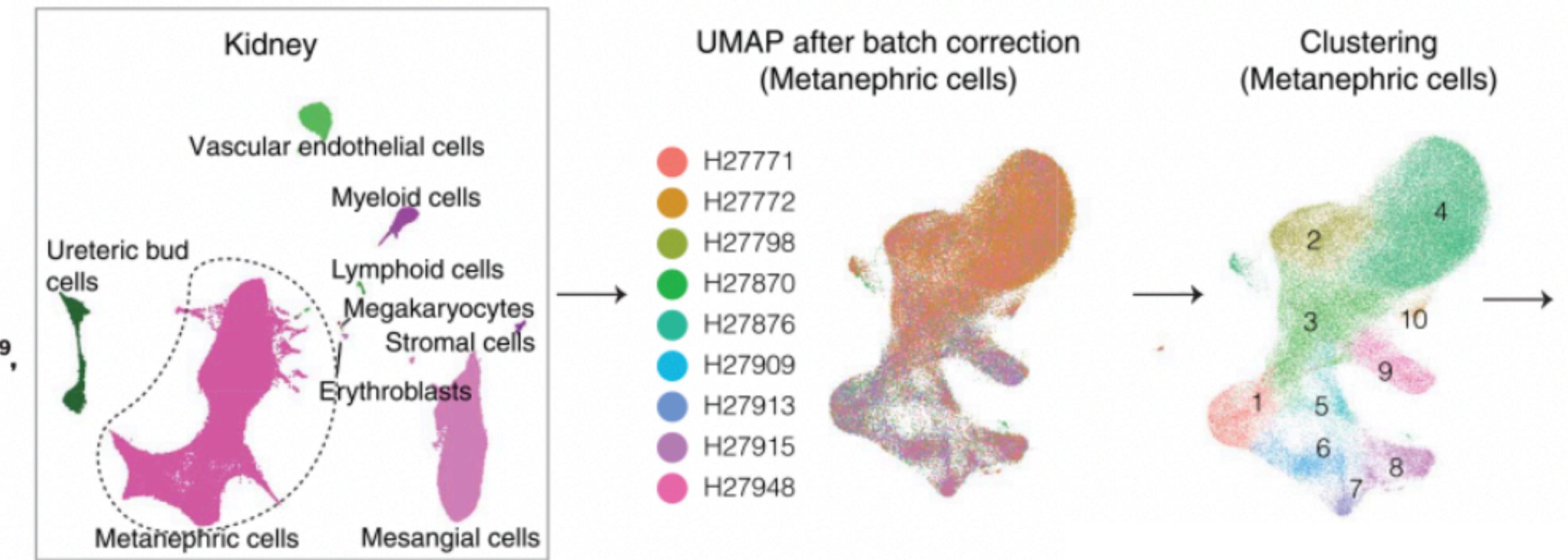
### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

## "Intradataset cross validation"

- Step 1: Cluster cells.
- Step 2: Treat clusters as truth.  
Do 5-fold cross validation with SVM.



# How can we validate the results of clustering?

## RESEARCH ARTICLE

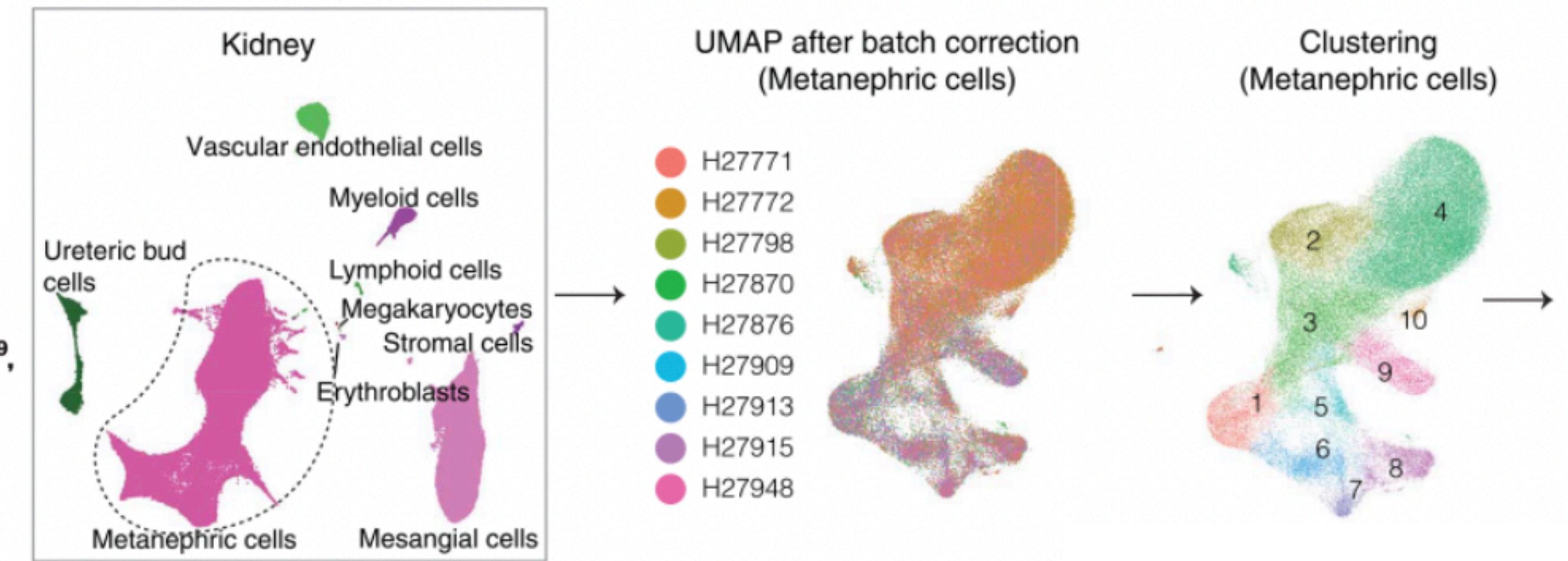
### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

## "Intradataset cross validation"

- **Step 1:** Cluster cells.
- **Step 2:** Treat clusters as truth.  
Do 5-fold cross validation with SVM.
- **Step 3:** Compare clusters to SVM predictions.



# How can we validate the results of clustering?

## RESEARCH ARTICLE

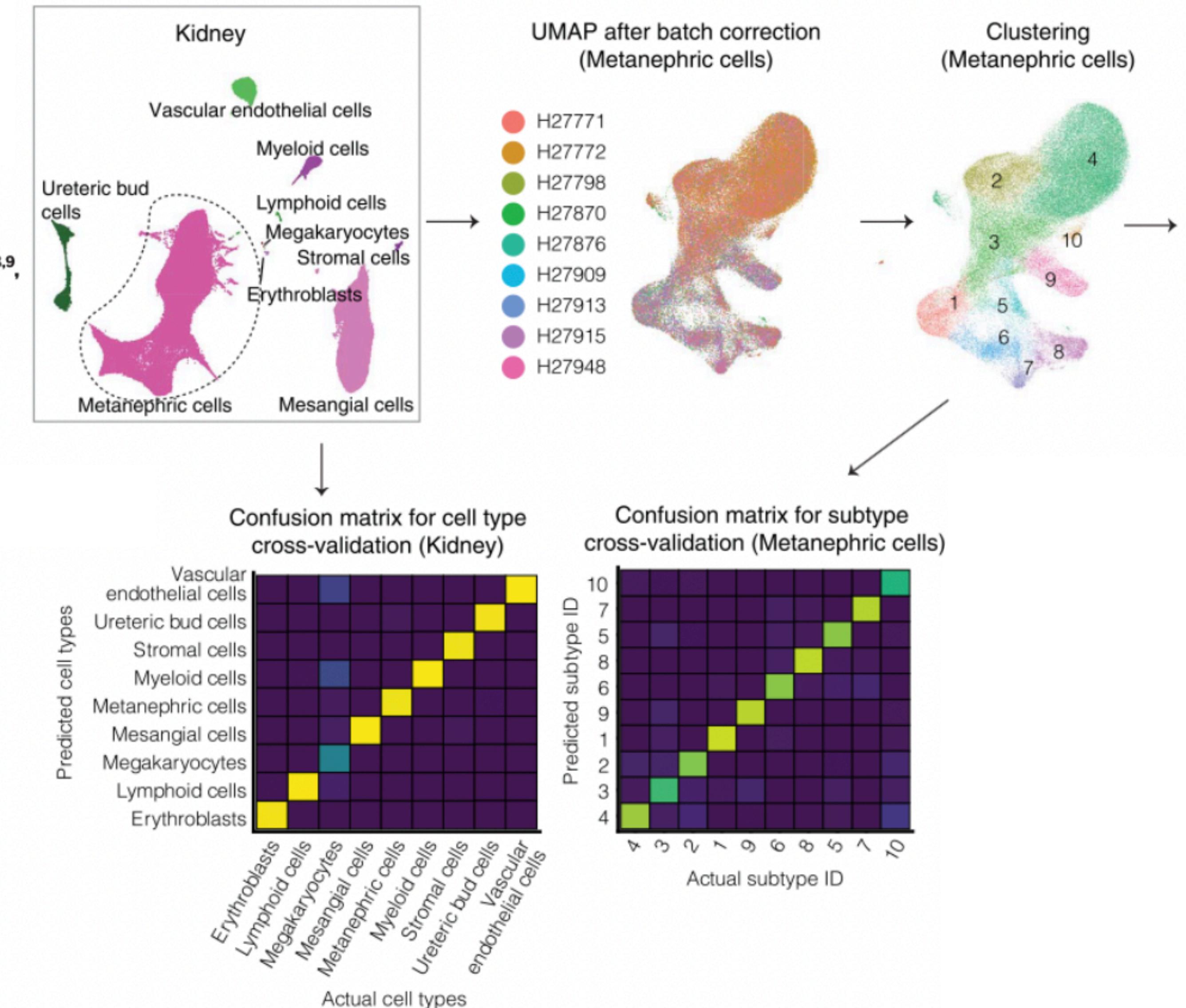
### HUMAN GENOMICS

## A human cell atlas of fetal gene expression

Junyue Cao<sup>1\*</sup>, Diana R. O'Day<sup>2</sup>, Hannah A. Pliner<sup>3</sup>, Paul D. Kingsley<sup>4</sup>, Mei Deng<sup>2</sup>, Riza M. Daza<sup>1</sup>, Michael A. Zager<sup>3,5</sup>, Kimberly A. Aldinger<sup>2,6</sup>, Ronnie Blecher-Gonen<sup>1</sup>, Fan Zhang<sup>7</sup>, Malte Spielmann<sup>8,9</sup>, James Palis<sup>4</sup>, Dan Doherty<sup>2,3,6</sup>, Frank J. Steemers<sup>7</sup>, Ian A. Glass<sup>2,3,6</sup>, Cole Trapnell<sup>1,3,10†</sup>, Jay Shendure<sup>1,3,10,11†</sup>

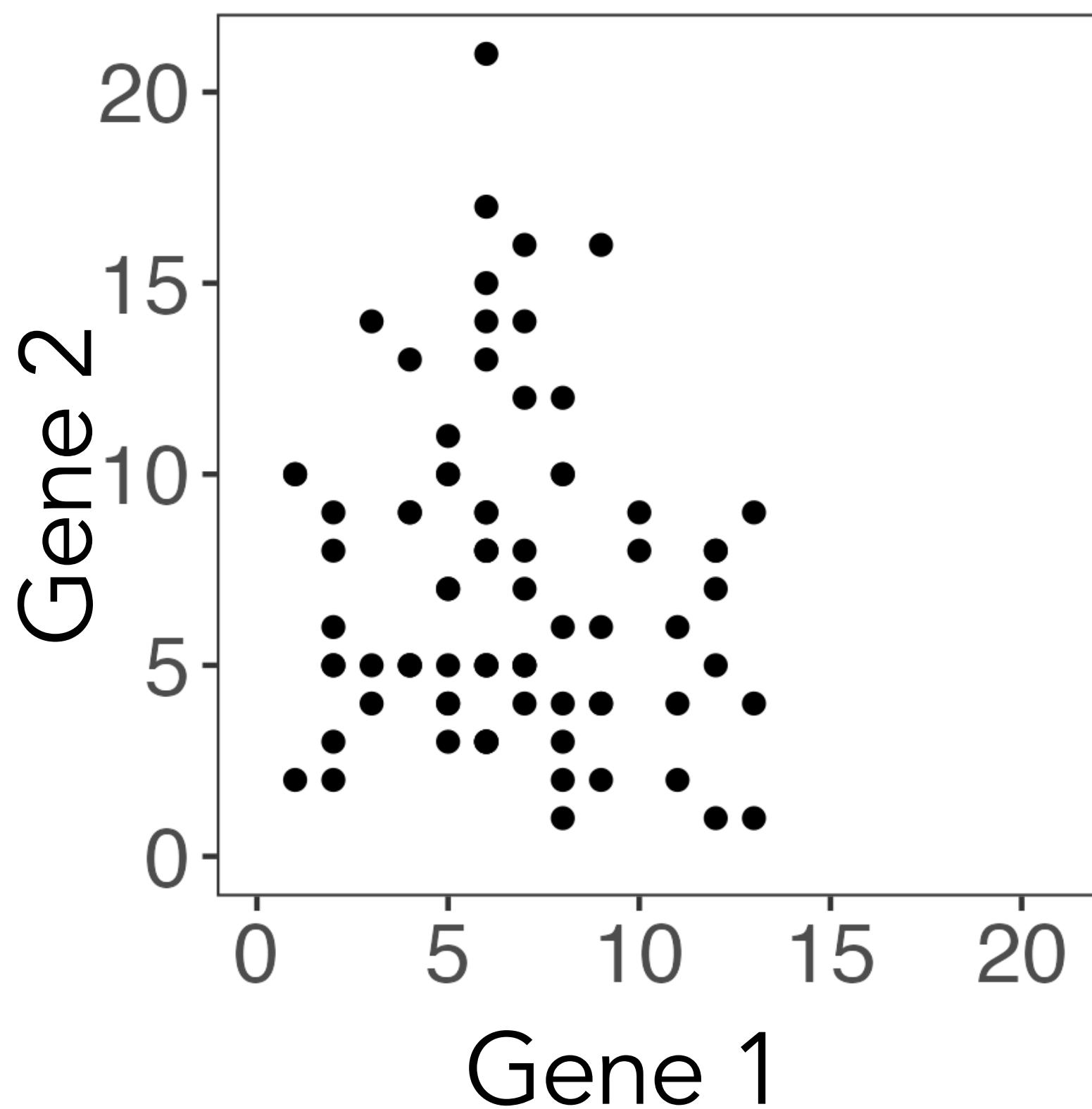
## "Intradataset cross validation"

- **Step 1:** Cluster cells.
- **Step 2:** Treat clusters as truth.  
Do 5-fold cross validation with SVM.
- **Step 3:** Compare clusters to SVM predictions.



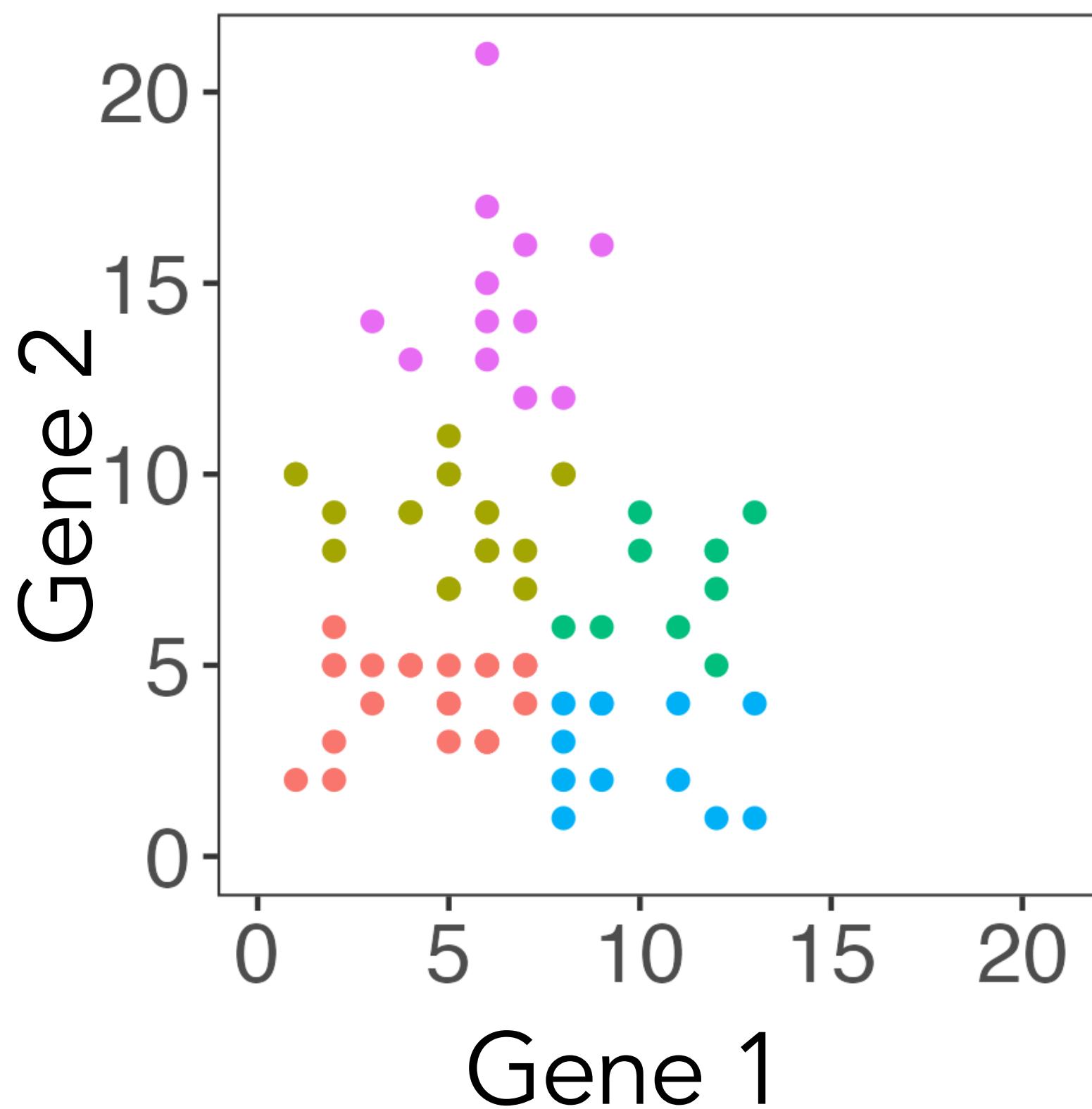
This cross validation procedure double dips

---

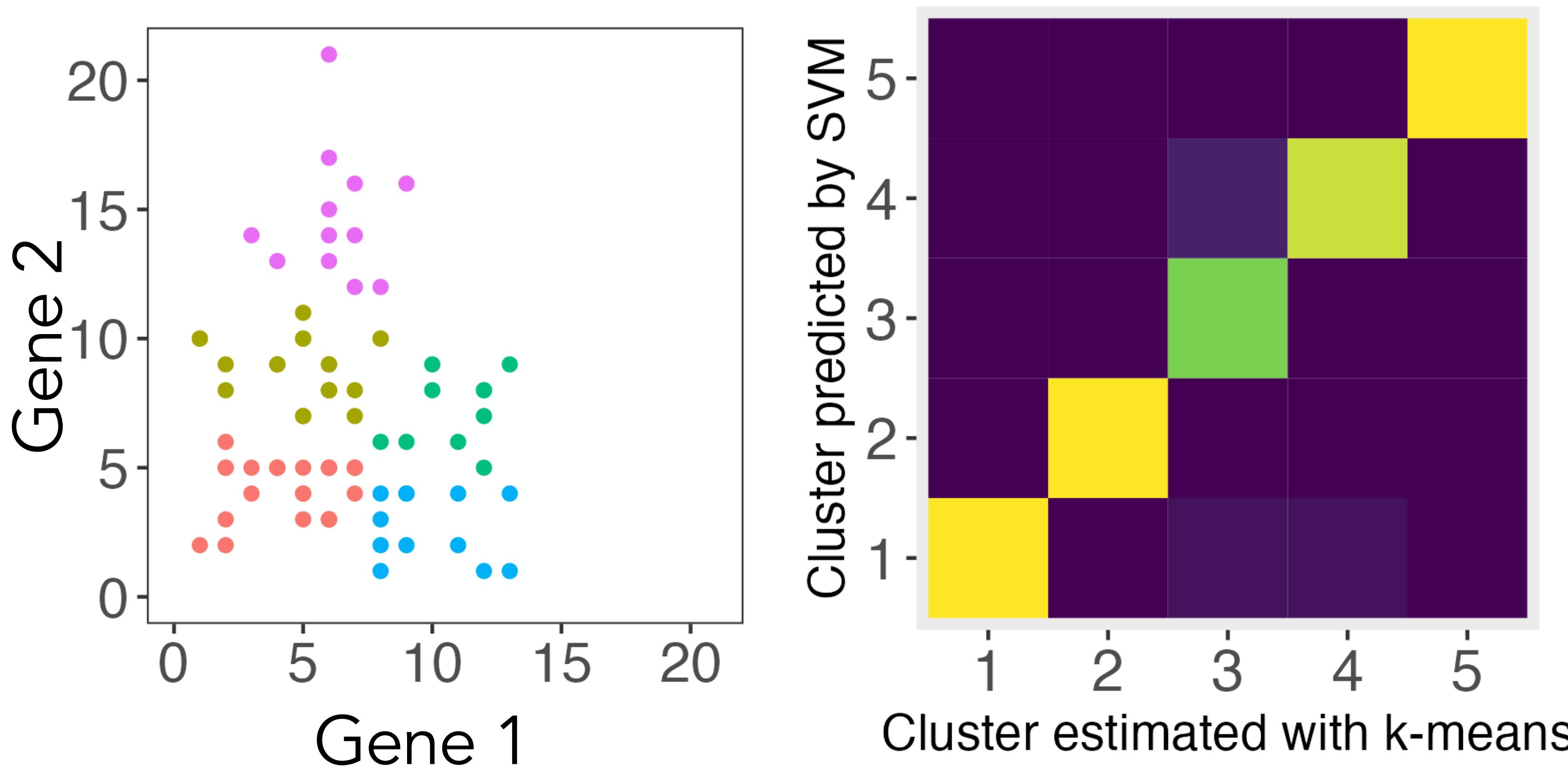


This cross validation procedure double dips

---



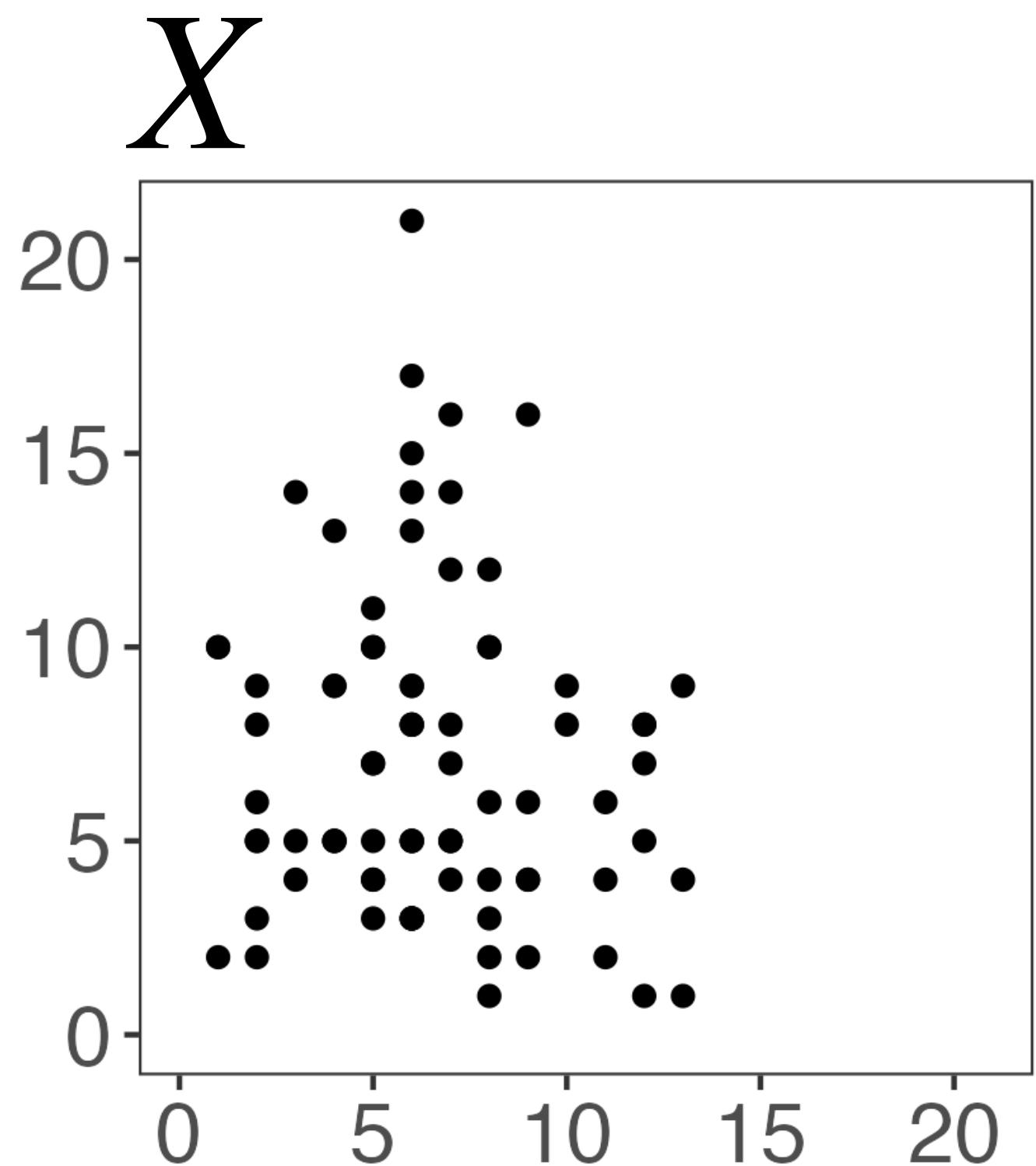
This cross validation procedure double dips



Classifier gets 96% accuracy to predict the five clusters, despite the fact that the five clusters are just random noise.

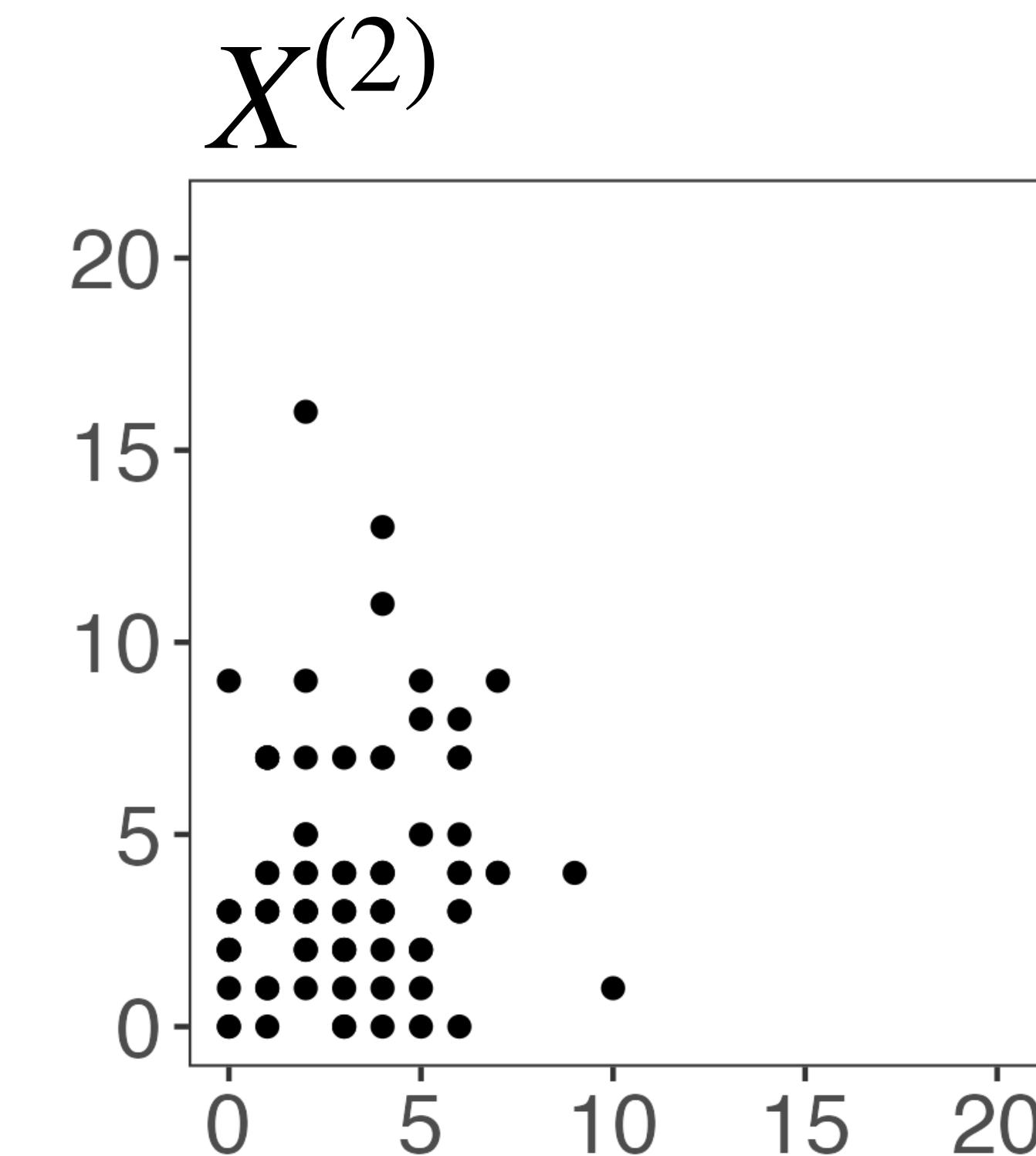
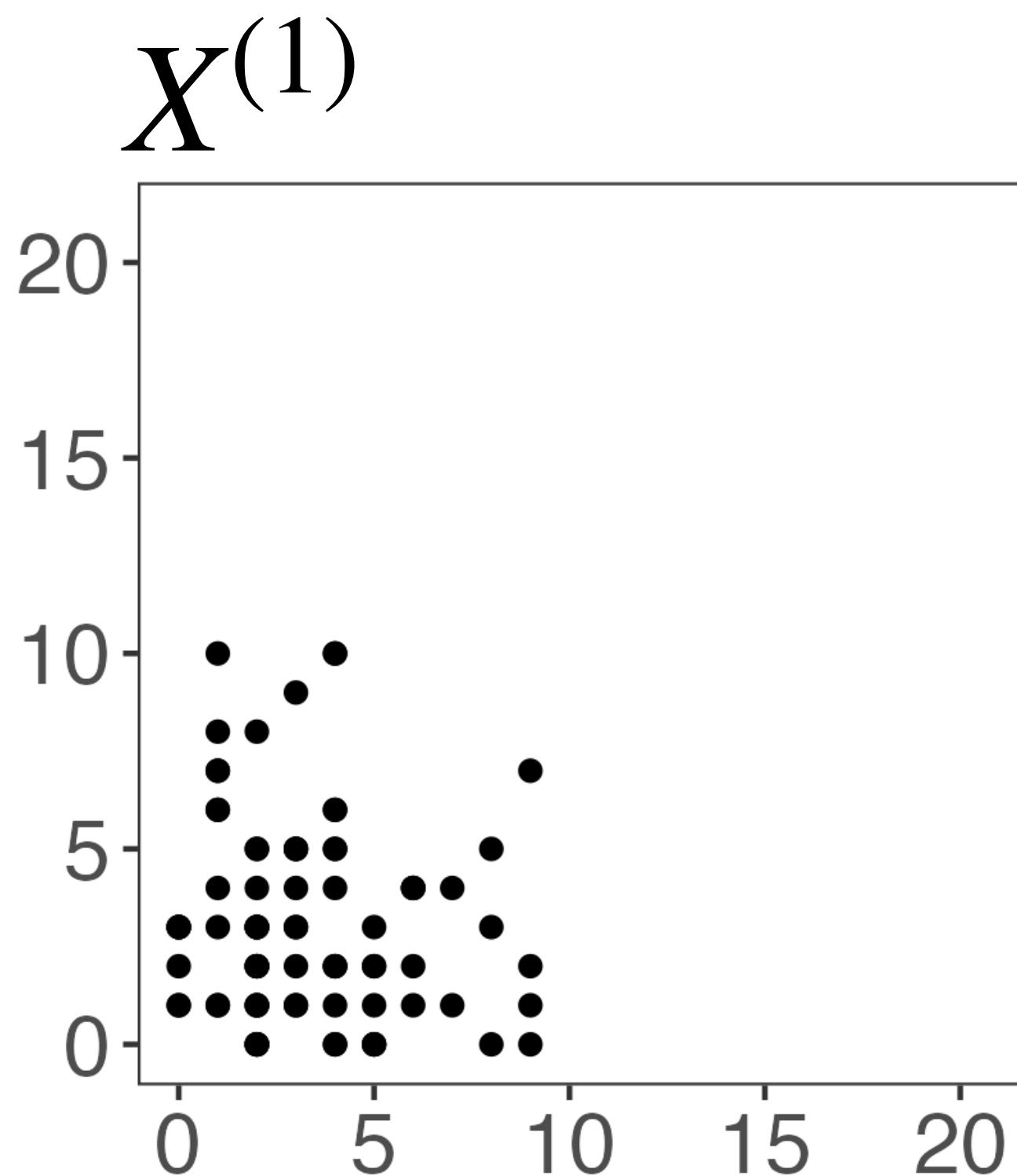
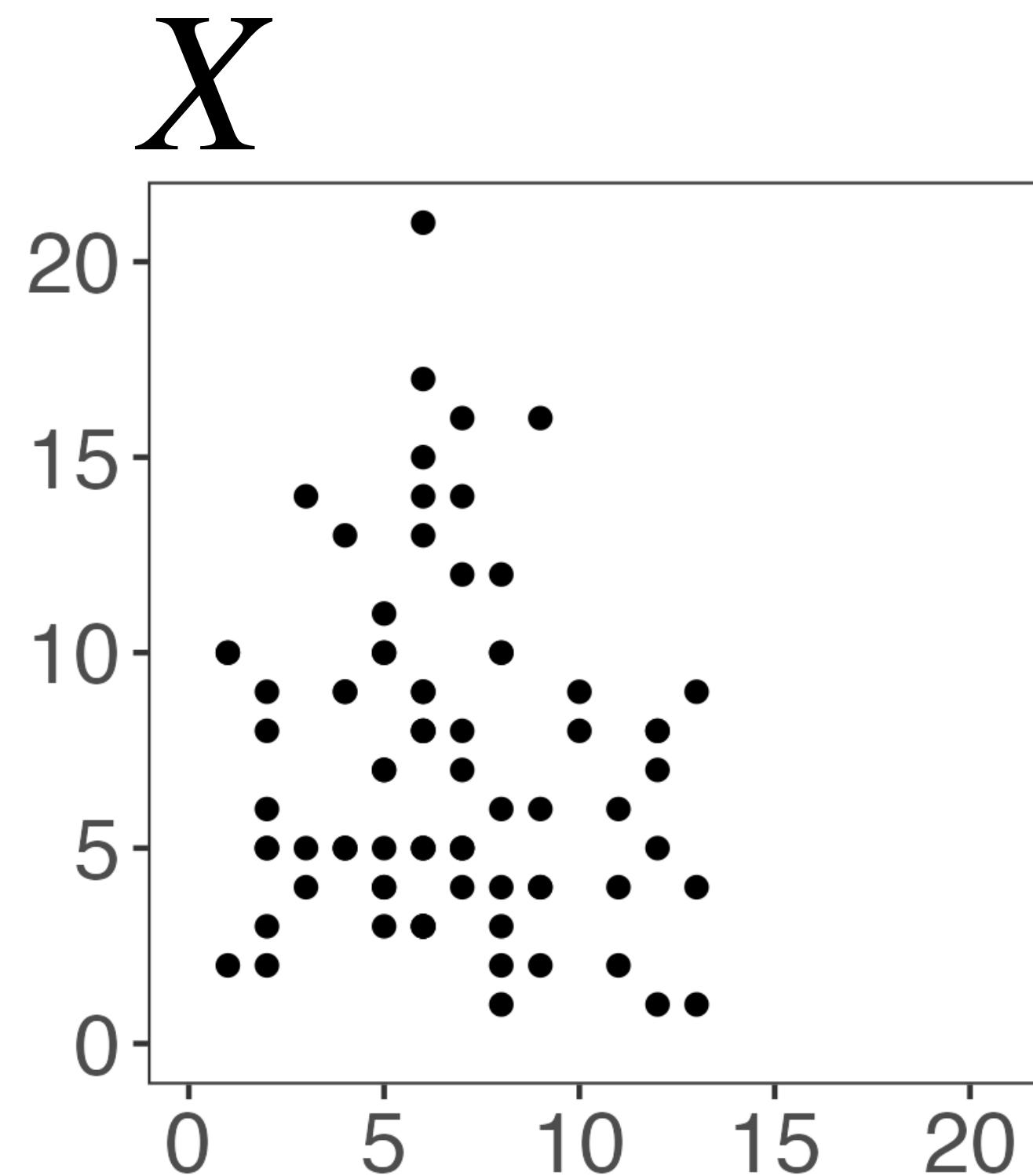
Data thinning provides a simple alternative

---



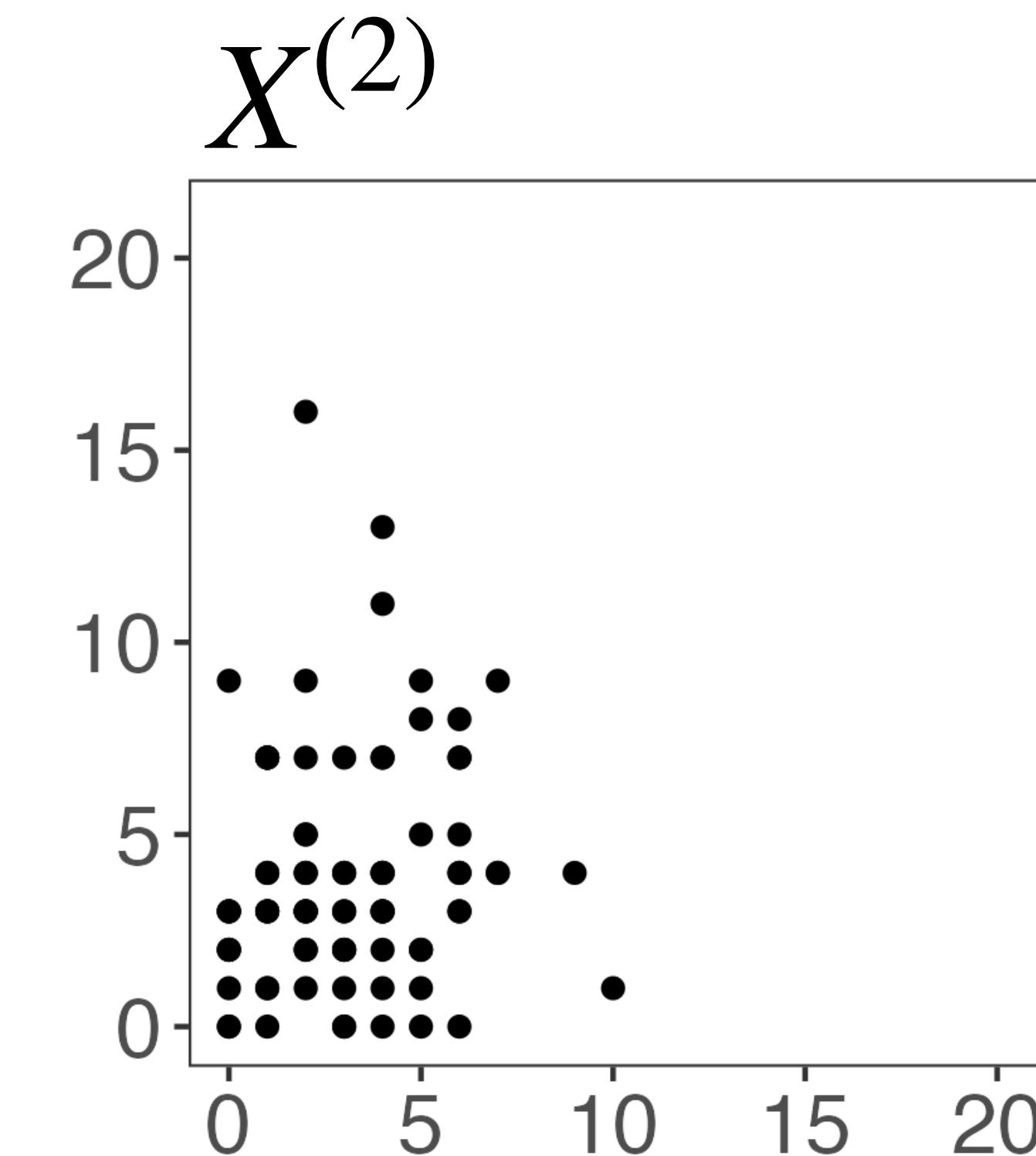
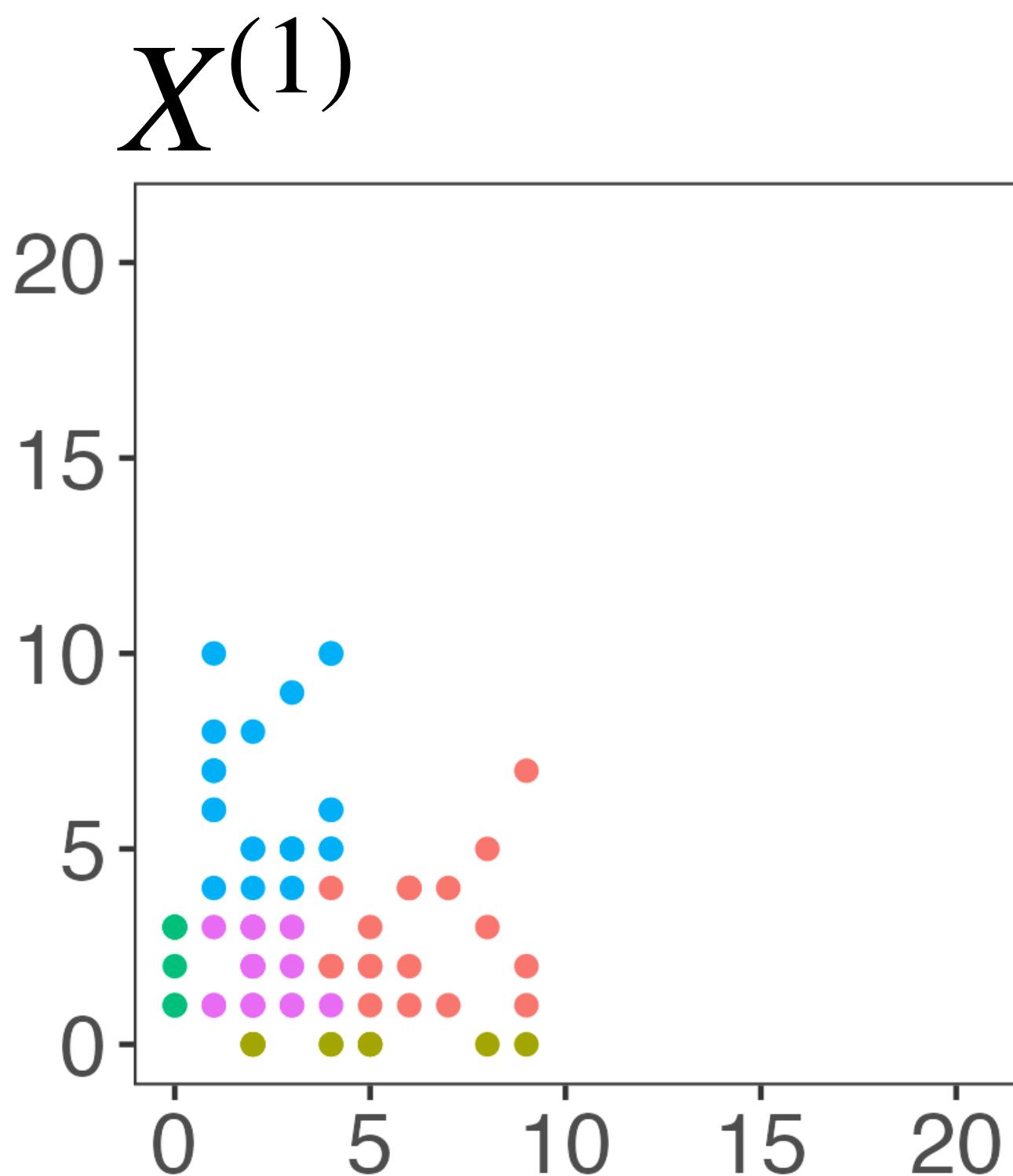
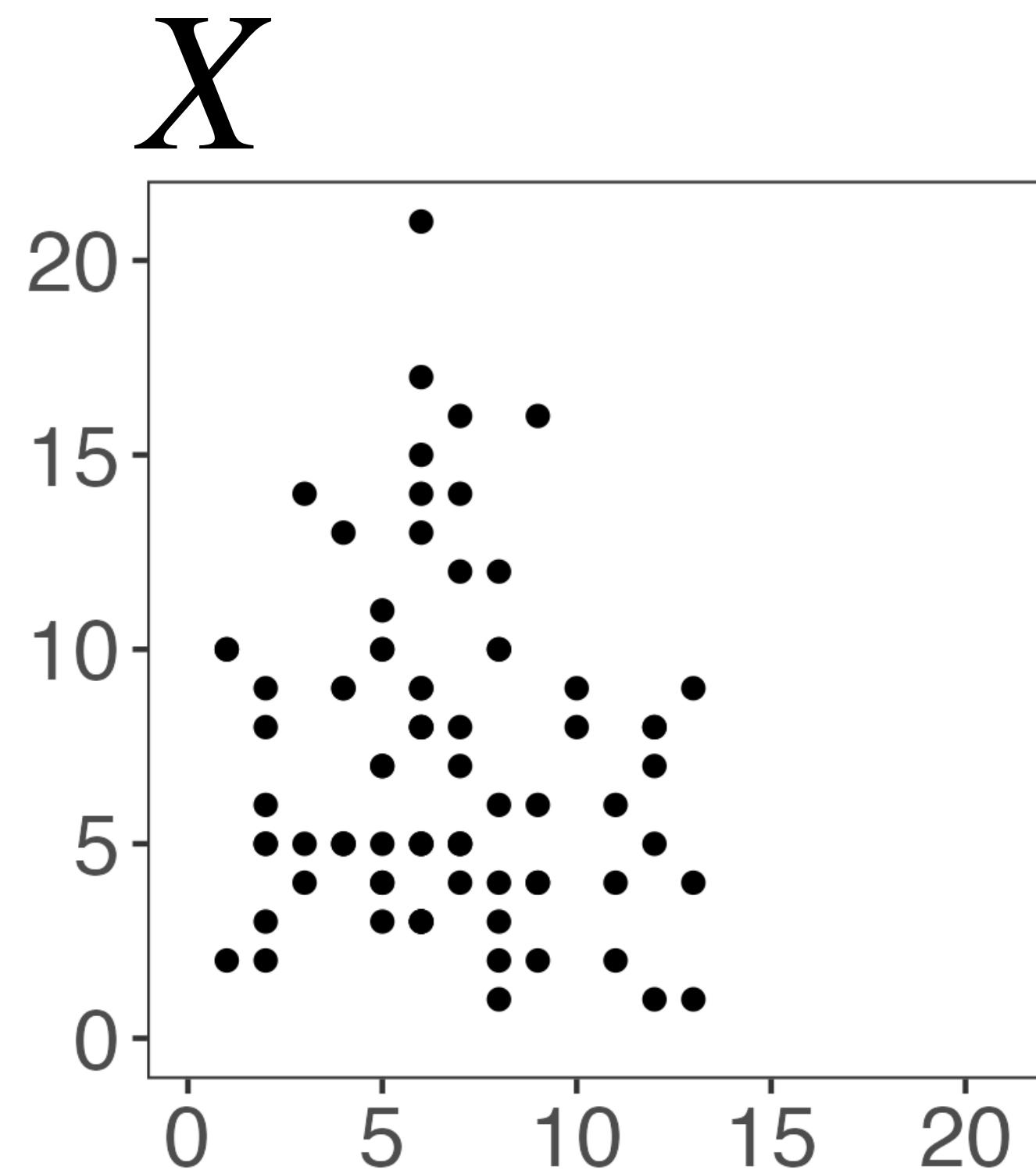
## Data thinning provides a simple alternative

---



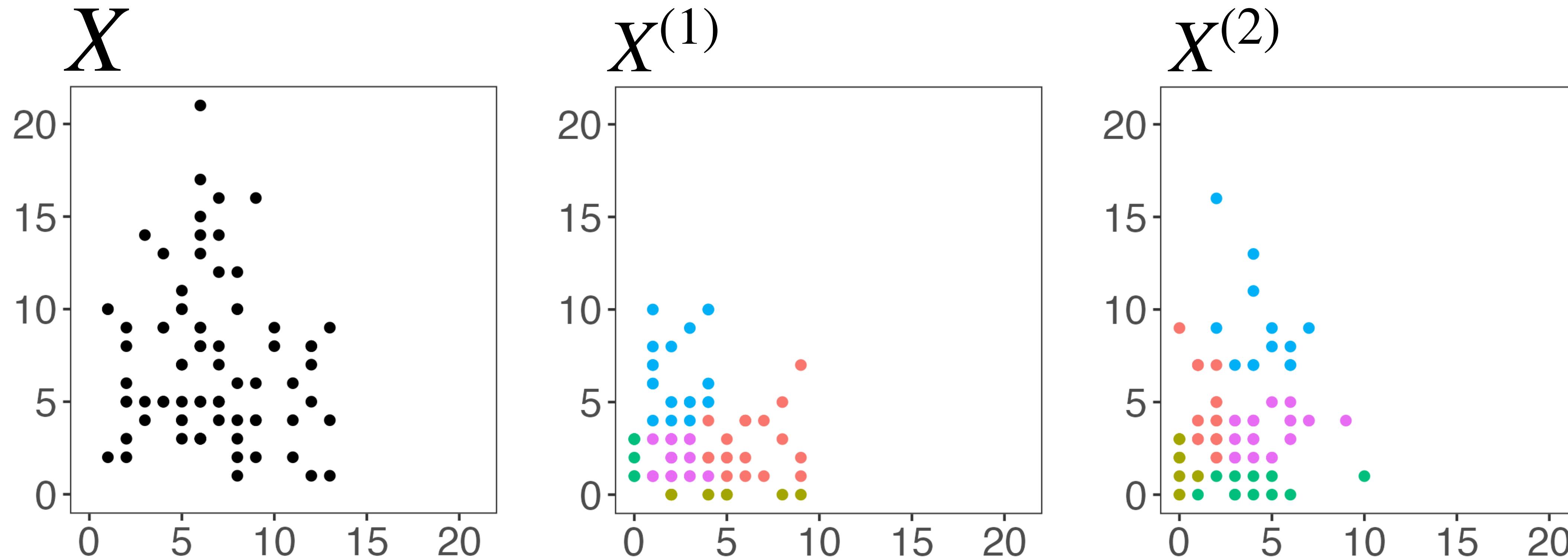
# Data thinning provides a simple alternative

---



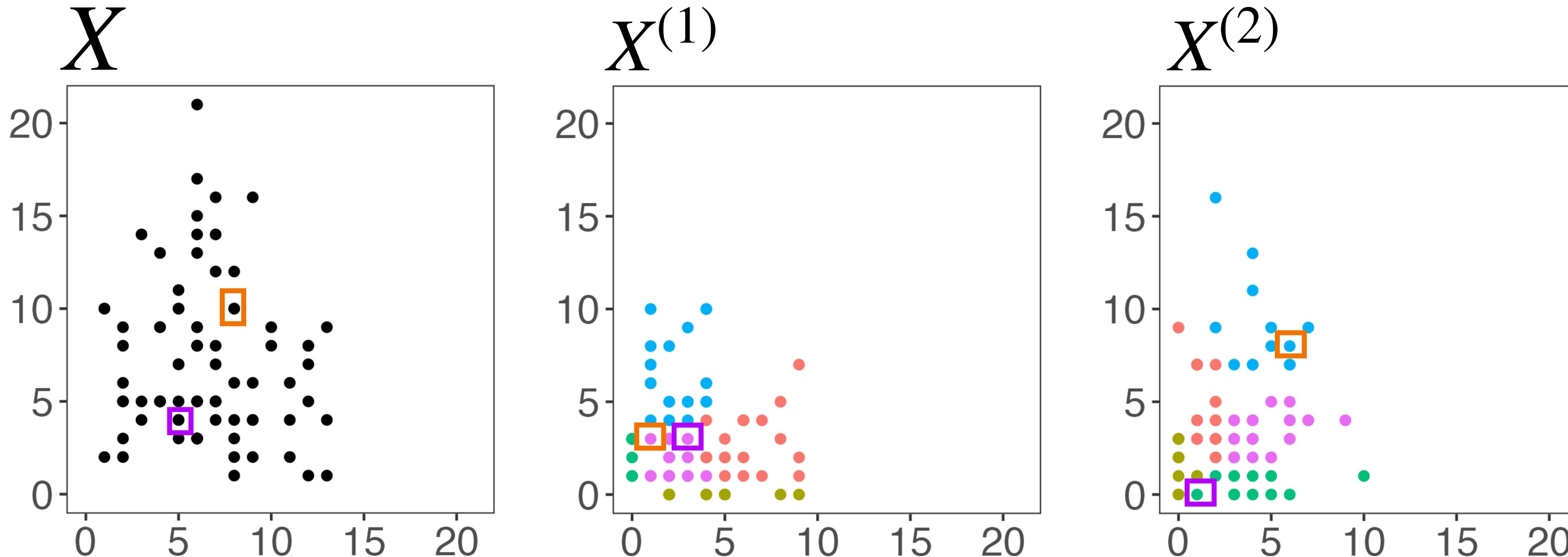
# Data thinning provides a simple alternative

---

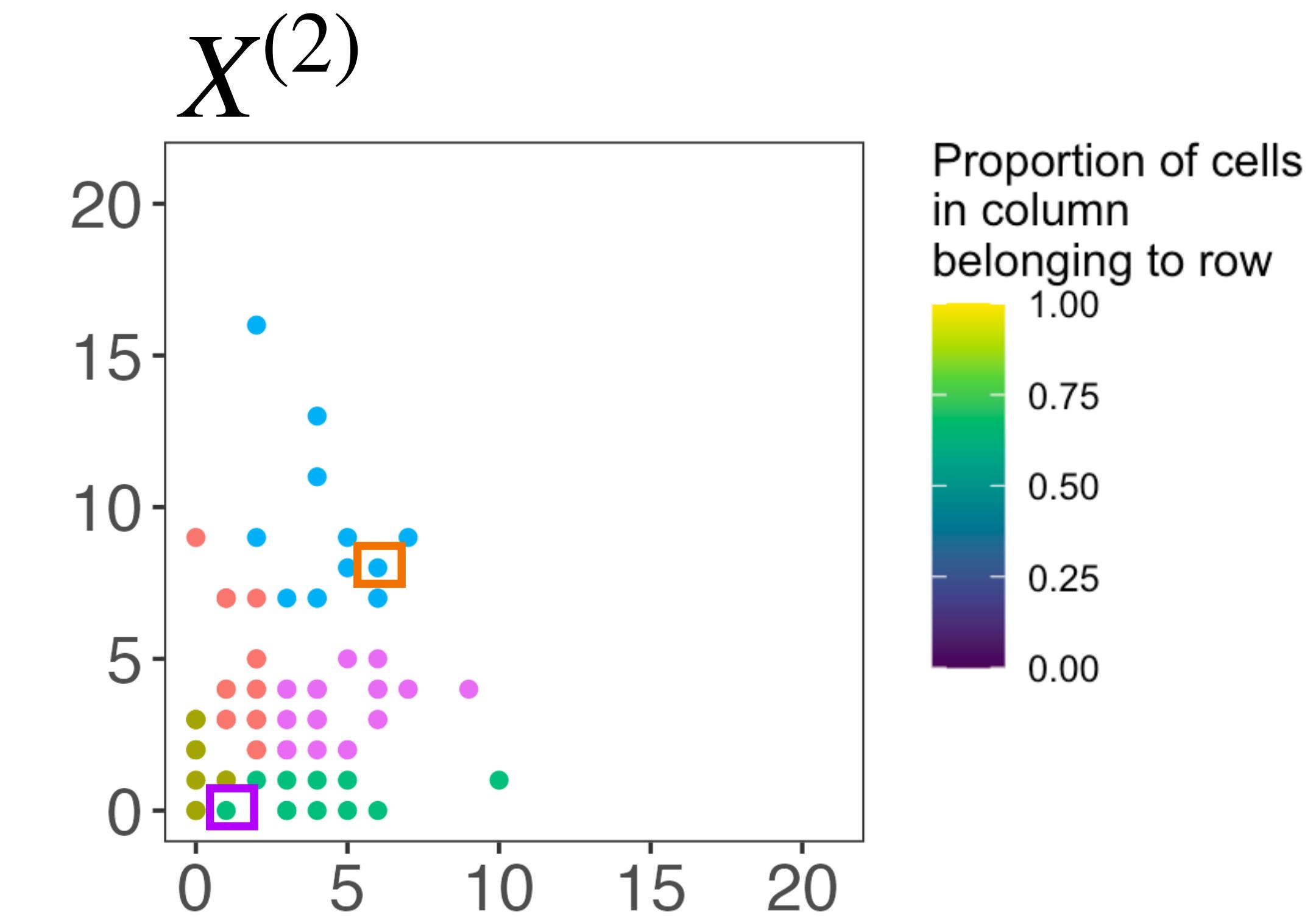
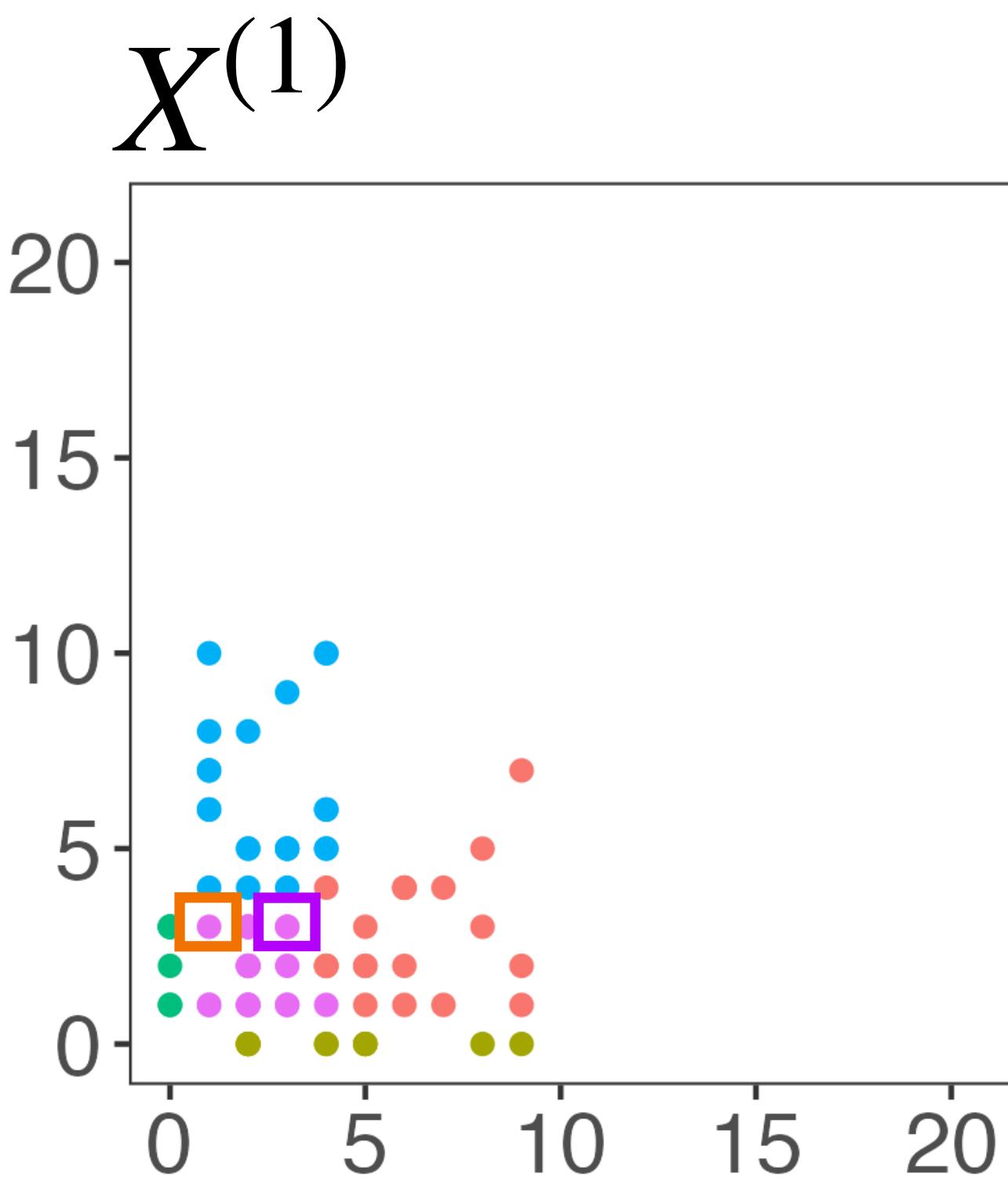
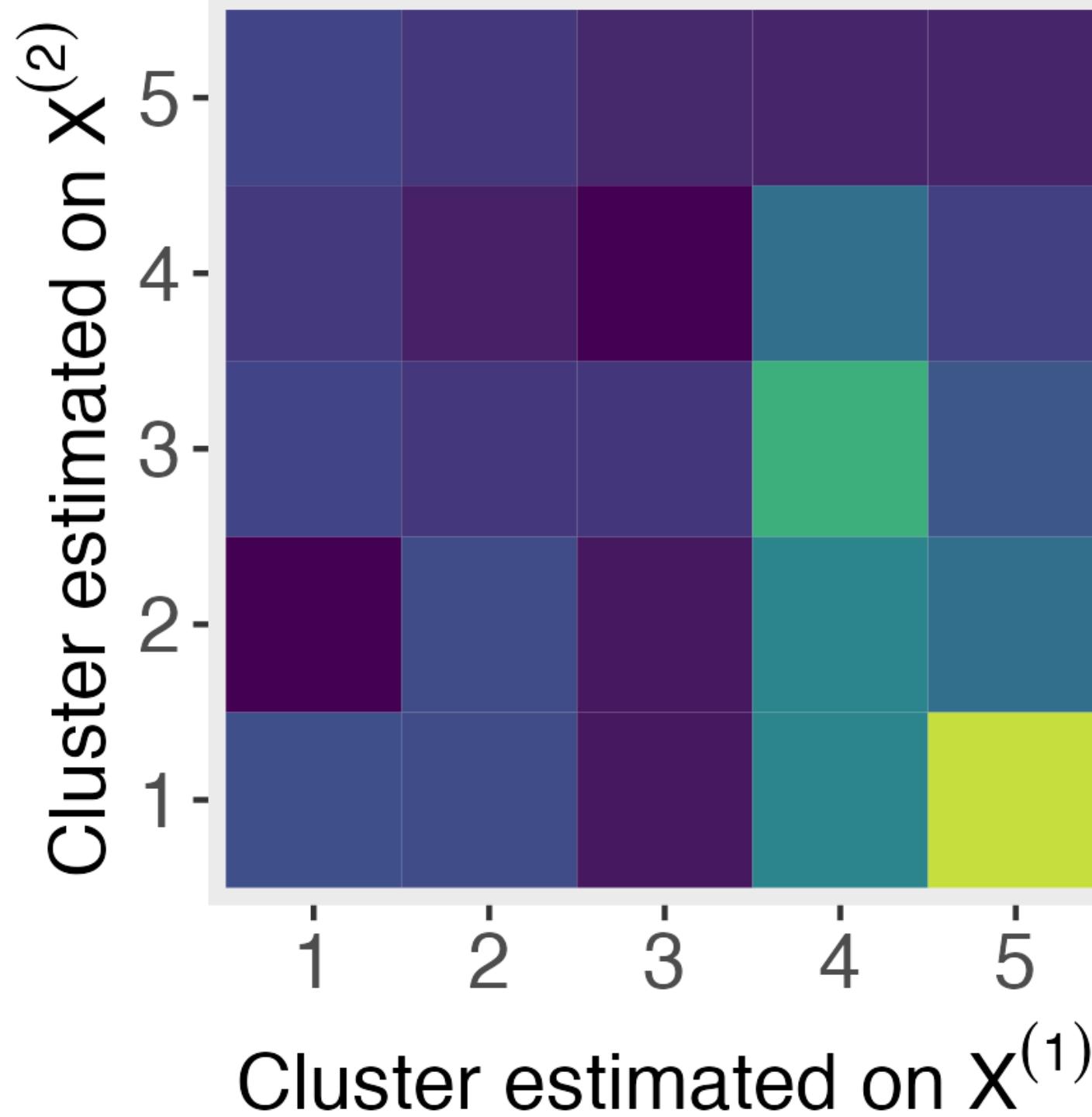


# Data thinning provides a simple alternative

---

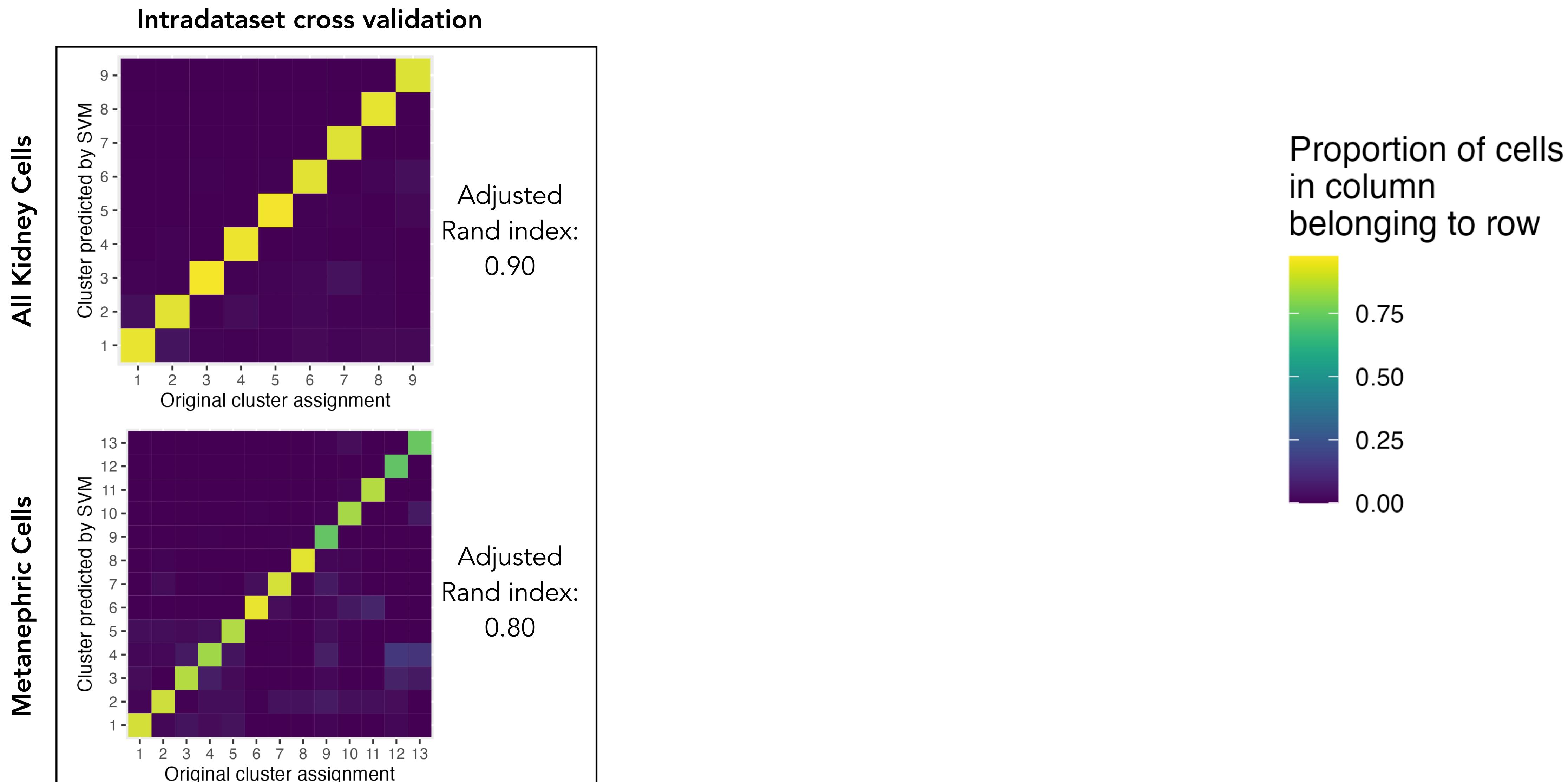


# Data thinning provides a simple alternative

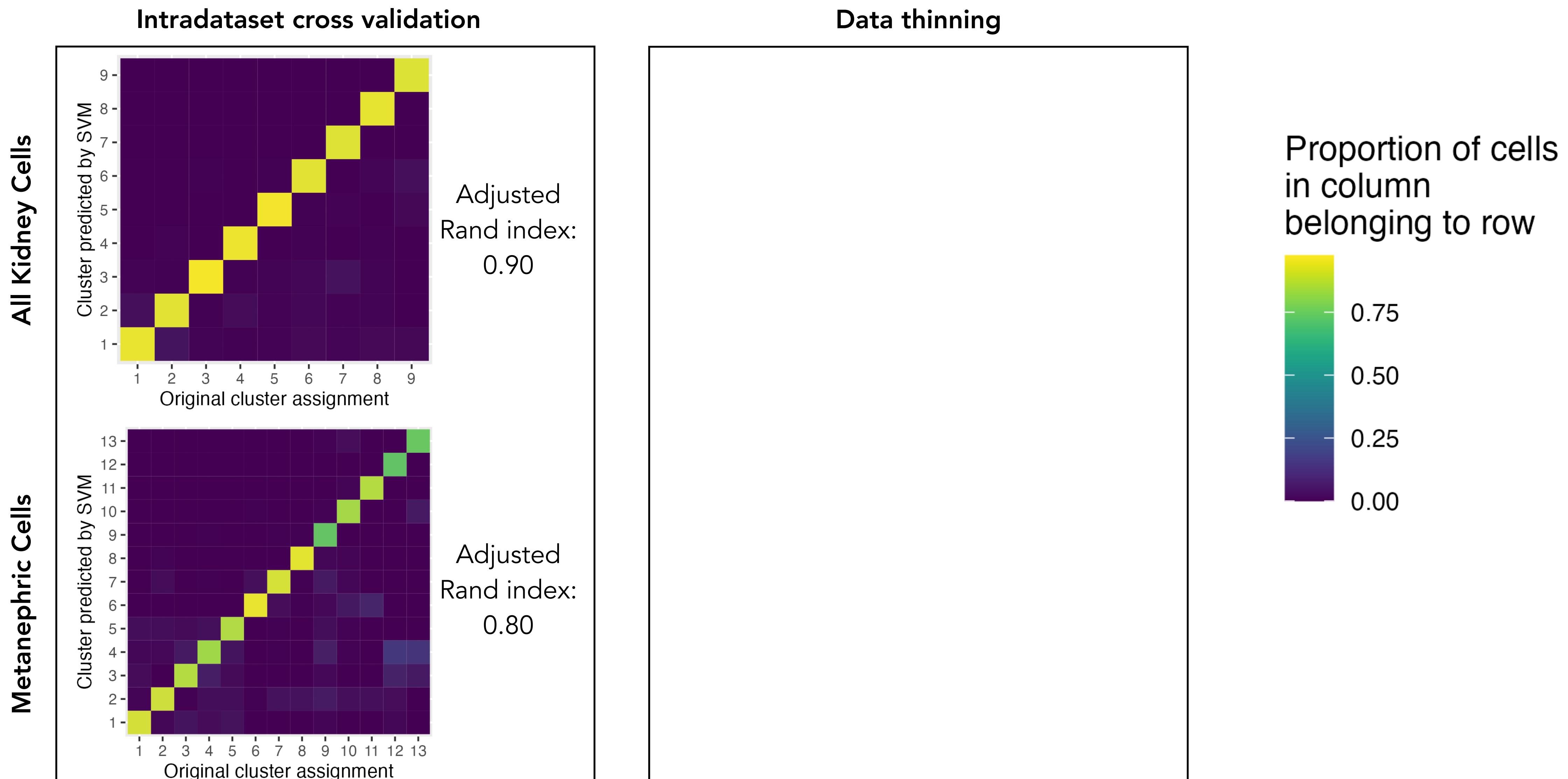


Adjusted Rand Index = 0.01

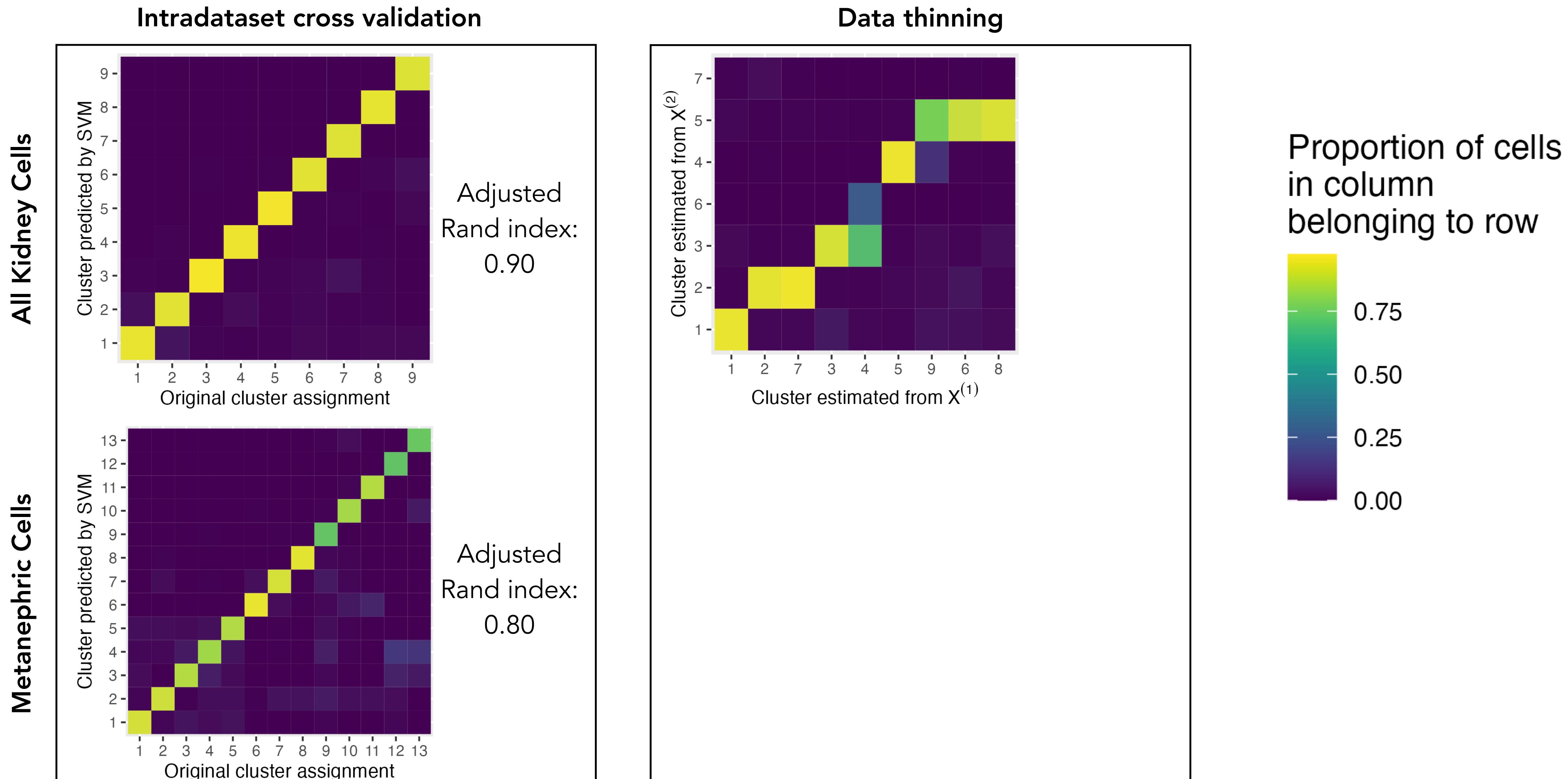
# Re-analysis of Kidney cell data from fetal cell atlas



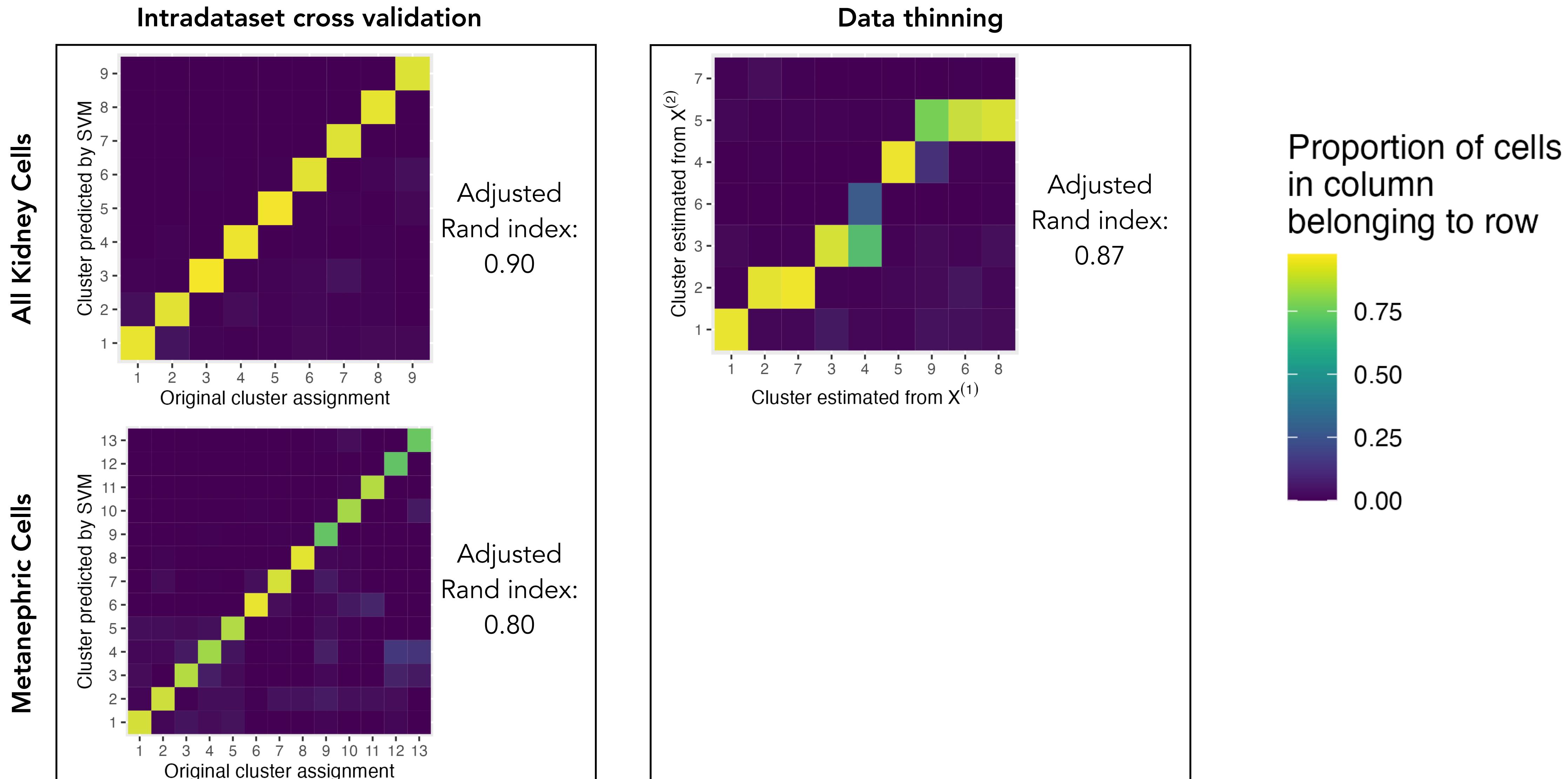
# Re-analysis of Kidney cell data from fetal cell atlas



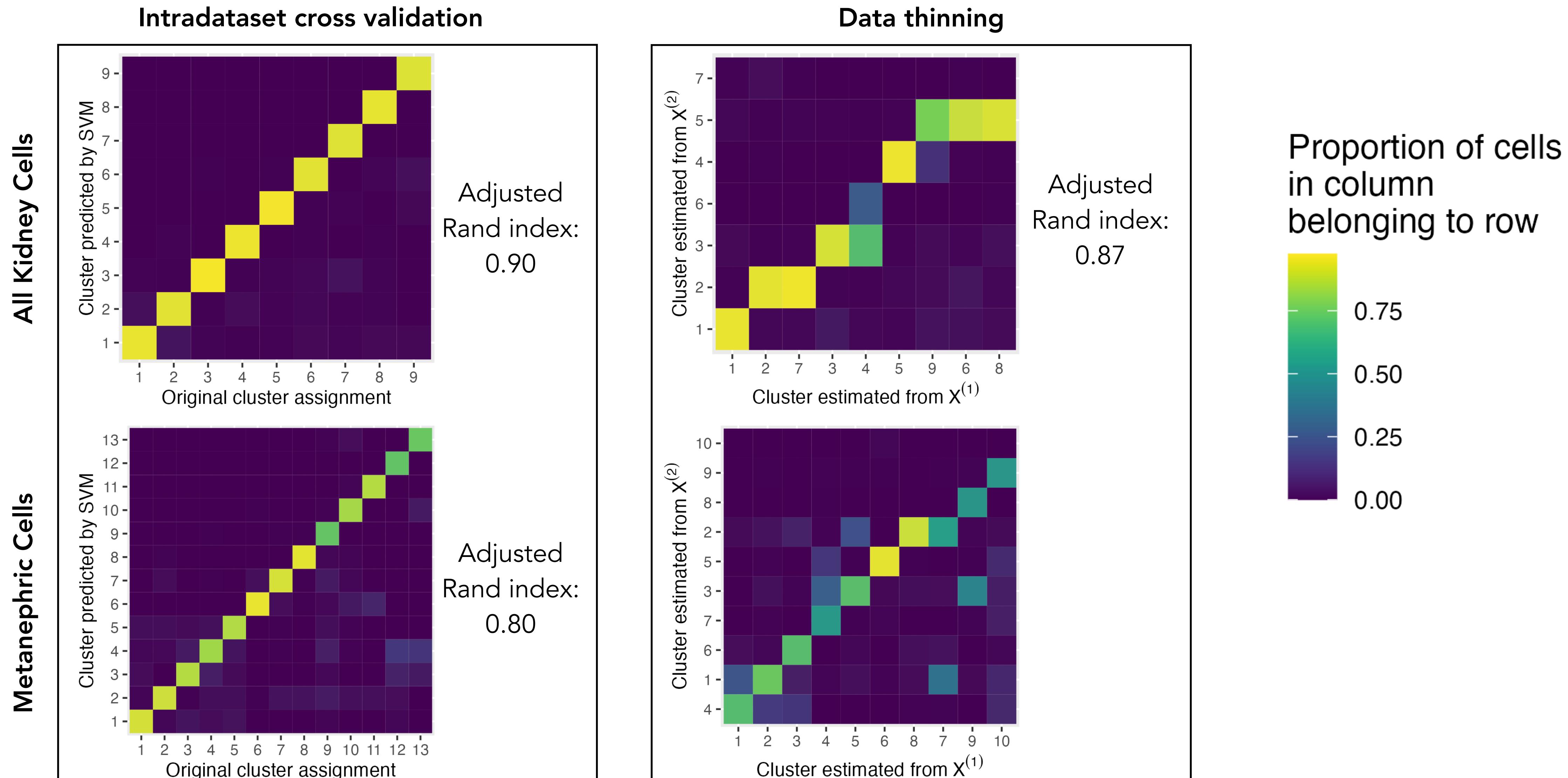
# Re-analysis of Kidney cell data from fetal cell atlas



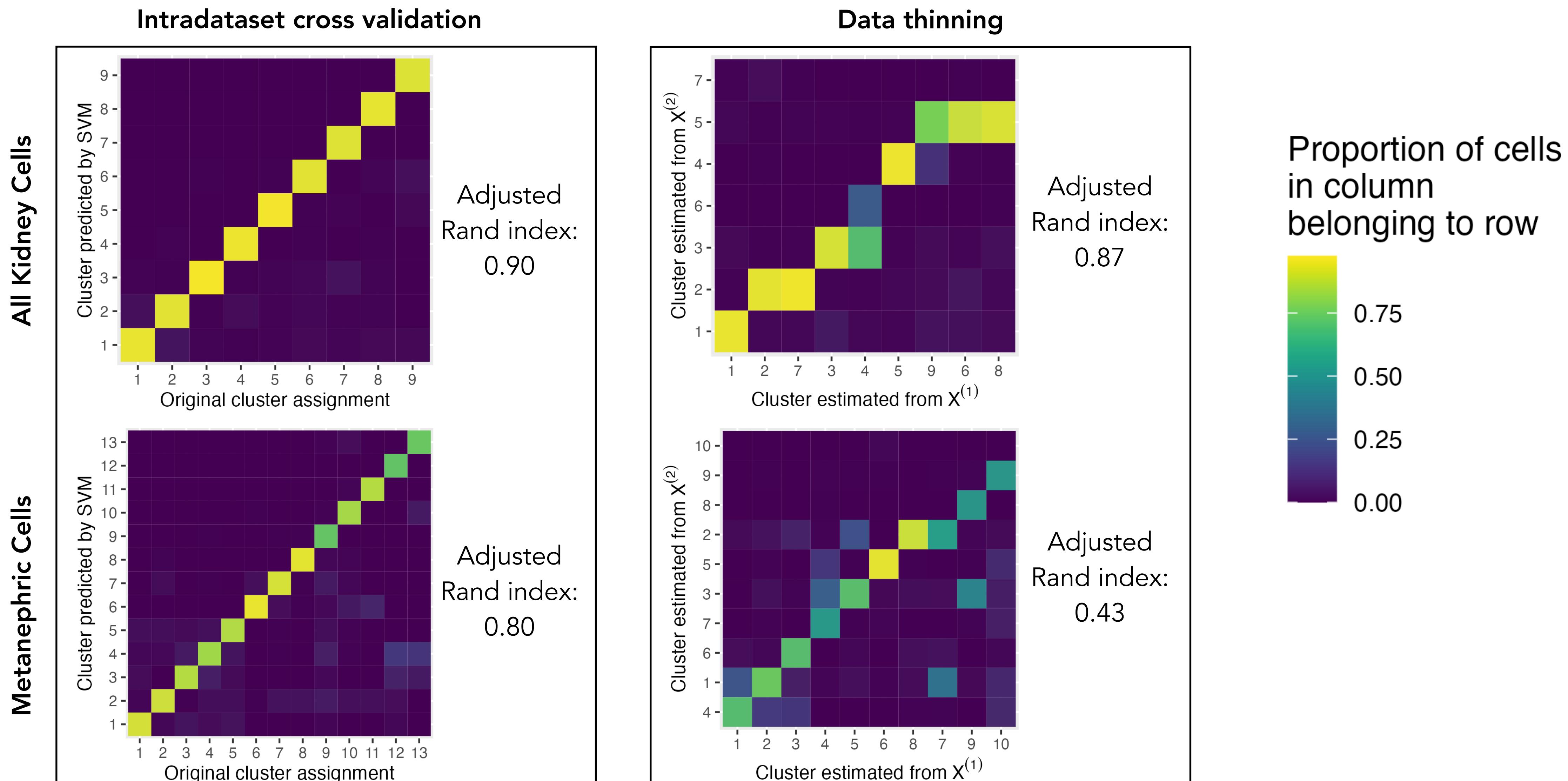
# Re-analysis of Kidney cell data from fetal cell atlas



# Re-analysis of Kidney cell data from fetal cell atlas



# Re-analysis of Kidney cell data from fetal cell atlas



## Outline

---

1. Motivation: settings where sample splitting doesn't work
2. Poisson thinning
3. Data thinning
4. Application to single-cell RNA sequencing data
5. **Ongoing work**

## Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.
2. Sample splitting.
3. Data thinning.

## Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.  
Requires a bespoke solution for every problem at hand.
2. Sample splitting.
3. Data thinning.

## Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.

Requires a bespoke solution for every problem at hand.

2. Sample splitting.

Super flexible! Requires only an iid assumption on your data.

3. Data thinning.

# Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.

Requires a bespoke solution for every problem at hand.

2. Sample splitting.

Super flexible! Requires only an iid assumption on your data.

Not an option in unsupervised settings; unsatisfying in other settings.

3. Data thinning.

# Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.

Requires a bespoke solution for every problem at hand.

2. Sample splitting.

Super flexible! Requires only an iid assumption on your data.

Not an option in unsupervised settings; unsatisfying in other settings.

3. Data thinning.

No bespoke solutions needed; works in supervised and unsupervised settings.

# Three ways to avoid double dipping

---

1. Specialized methods, such as selective inference.

Requires a bespoke solution for every problem at hand.

2. Sample splitting.

Super flexible! Requires only an iid assumption on your data.

Not an option in unsupervised settings; unsatisfying in other settings.

3. Data thinning.

No bespoke solutions needed; works in supervised and unsupervised settings.

Requires distributional assumptions and knowledge of nuisance parameters.

# Three ways to avoid double dipping

---

## 1. Specialized methods, such as selective inference.

Requires a bespoke solution for every problem at hand.

## 2. Sample splitting.

Super flexible! Requires only an iid assumption on your data.

Not an option in unsupervised settings; unsatisfying in other settings.

## 3. Data thinning.

No bespoke solutions needed; works in supervised and unsupervised settings.

Requires distributional assumptions and knowledge of nuisance parameters.

Limited to convolution-closed distributions?

## Revisiting the goals of data thinning

---

**Goal:** split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

## Revisiting the goals of data thinning

---

**Goal:** split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

In our previous recipe:

- (3)  $X = X^{(1)} + X^{(2)}$ .

## Revisiting the goals of data thinning

---

**Goal:** split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

In our previous recipe:

~~(3)  $X = X^{(1)} + X^{(2)}$ .~~

## Revisiting the goals of data thinning

---

**Goal:** split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- (1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

In our previous recipe:

~~(3)  $X = X^{(1)} + X^{(2)}$ .~~ (3)  $X = T(X^{(1)}, X^{(2)})$ .

## Revisiting the goals of data thinning

---

**Goal:** split a single observation  $X$  into  $X^{(1)}$  and  $X^{(2)}$  such that:

- ~~(1)  $X^{(1)}$  and  $X^{(2)}$  have the same distribution as  $X$ , up to a parameter scaling.~~
- (2)  $X^{(1)} \perp\!\!\!\perp X^{(2)}$ .

In our previous recipe:

- ~~(3)  $X = X^{(1)} + X^{(2)}$ .~~ (3)  $X = T(X^{(1)}, X^{(2)})$ .

# Generalized thinning with non-additive decompositions

---

# Generalized thinning with non-additive decompositions

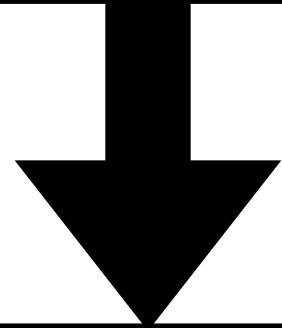
---

We observe realization  $x$  from  $X \sim P_\theta$ .

# Generalized thinning with non-additive decompositions

---

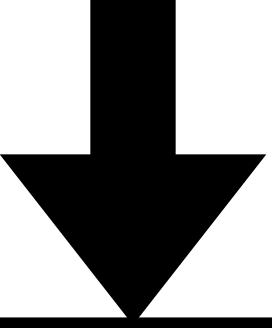
We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim P_\theta$ .

# Generalized thinning with non-additive decompositions

We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .

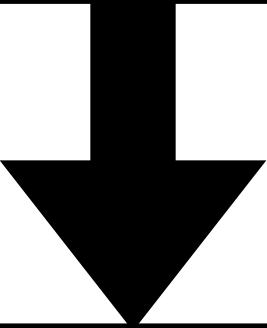


Can we work backwards to recover  
 $x'$  and  $x''$ ?

We observe realization  $x$  from  $X \sim P_\theta$ .

# Generalized thinning with non-additive decompositions

We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .



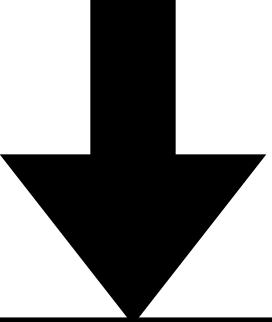
We observe realization  $x$  from  $X \sim P_\theta$ .

Can we work backwards to recover  
 $x'$  and  $x''$ ?

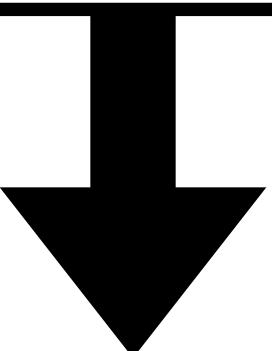
Let  $G_{x,\theta}$  be the conditional distribution of  
 $(X', X'') \mid X = x$ .

# Generalized thinning with non-additive decompositions

We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim P_\theta$ .



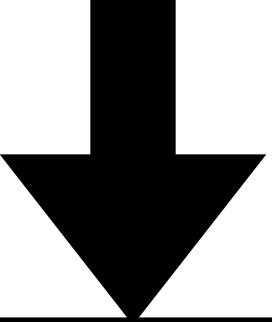
Draw  $(X^{(1)}, X^{(2)})$  from  $G_{x,\theta}$ .

Can we work backwards to recover  $x'$  and  $x''$ ?

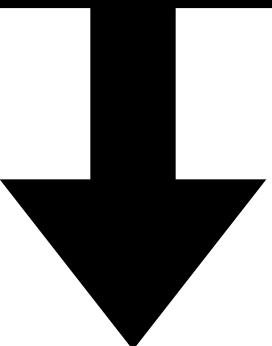
Let  $G_{x,\theta}$  be the conditional distribution of  $(X', X'') \mid X = x$ .

# Generalized thinning with non-additive decompositions

We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim P_\theta$ .



Draw  $(X^{(1)}, X^{(2)})$  from  $G_{x,\theta}$ .

Can we work backwards to recover  $x'$  and  $x''$ ?

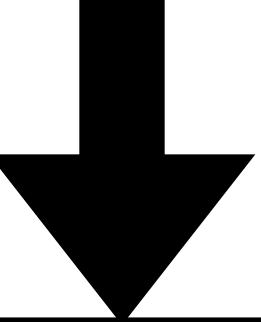
Let  $G_{x,\theta}$  be the conditional distribution of  $(X', X'') \mid X = x$ .

**Theorem:**

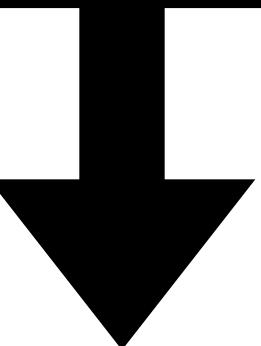
$$X^{(1)} \sim Q_\theta^1, \quad X^{(2)} \sim Q_\theta^2, \quad X^{(1)} \perp\!\!\!\perp X^{(2)}.$$

# Generalized thinning with non-additive decompositions

We know  $x$  could have arisen as  $T(x', x'')$ , where  
 $X' \sim Q_\theta^1$ ,  $X'' \sim Q_\theta^2$ ,  $X' \perp\!\!\!\perp X''$ .



We observe realization  $x$  from  $X \sim P_\theta$ .



Draw  $(X^{(1)}, X^{(2)})$  from  $G_{x,\theta}$ .

**Theorem:**

$$X^{(1)} \sim Q_\theta^1, \quad X^{(2)} \sim Q_\theta^2, \quad X^{(1)} \perp\!\!\!\perp X^{(2)}.$$

Can we work backwards to recover  $x'$  and  $x''$ ?

Let  $G_{x,\theta}$  be the conditional distribution of  $(X', X'') \mid X = x$ .

**Key idea:** If  $X = T(X', X'')$  is sufficient for  $\theta$  in the joint of  $(X', X'')$ , then  $G_{x,\theta}$  does not depend on  $\theta$ .

# The list of distributions we can thin is extensive

Family	Distribution $P_\theta$ , where $X \sim P_\theta$ .	Distribution $Q_\theta^{(k)}$ where $X^{(k)} \stackrel{ind.}{\sim} Q_\theta^{(k)}$ .	Sufficient statistic $T$ (sufficient for $\theta$ )
Natural exponential family (in parameter $\theta$ )	$N(\theta, \sigma^2)$	$N(\epsilon_k \theta, \epsilon_k \sigma^2)$	
	Poisson( $\theta$ )	Poisson( $\epsilon_k \theta$ )	
	NegBin( $r, \theta$ )	NegBin( $\epsilon_k r, \theta$ )	
	Binomial( $r, \theta$ )	Binomial( $\epsilon_k r, \theta$ )	$\sum_{k=1}^K X^{(k)}$
	Gamma( $\alpha, \theta$ )	Gamma( $\epsilon_k \alpha, \theta$ )	
	$N_p(\boldsymbol{\theta}, \Sigma)$	$N_p(\epsilon_k \boldsymbol{\theta}, \epsilon_k \Sigma)$	
	Multinomial $_p(r, \boldsymbol{\theta})$	Multinomial $_p(\epsilon_k r, \boldsymbol{\theta})$	
General exponential family (in parameter $\theta$ )	Gamma( $K/2, \theta$ )	$N(0, \frac{1}{2\theta})$	$\sum_{k=1}^K (X^{(k)})^2$
	Gamma( $K, \theta$ )	Weibull( $\theta^{-\frac{1}{\nu}}, \nu$ )	$\sum_{k=1}^K (X^{(k)})^\nu$
	Beta( $\theta, \beta$ )	Beta( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
	Beta( $\alpha, \theta$ )	Beta( $\frac{1}{K}\alpha, \frac{1}{K}\theta + \frac{k-1}{K}$ )	$(\prod_{k=1}^K (1 - X^{(k)}))^{1/K}$
	Gamma( $\theta, \beta$ )	Gamma( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
Truncated support family	Weibull( $\theta, \nu$ )	Gamma( $\frac{1}{K}, \theta^{-\nu}$ )	$(\sum_{k=1}^K X^{(k)})^{1/\nu}$
	Pareto( $\nu, \theta$ )	Gamma( $\frac{1}{K}, \theta$ )	$\nu \times \text{Exp}(\sum_{k=1}^K X^{(k)})$
	$N(0, \theta)$	Gamma( $\frac{1}{2K}, \frac{1}{2\theta}$ )	$X^2 = \sum_{k=1}^K X^{(k)}$
Non-parametric	$N_K(\theta_1 \mathbf{1}_K, \theta_2 I_K)$	$N(\theta_1, \theta_2)$	sample mean and variance
	Unif( $0, \theta$ )	$\theta \cdot \text{Beta}(\frac{1}{K}, 1)$	$\max(X^{(1)}, \dots, X^{(K)})$
	$\theta \cdot \text{Beta}(\alpha, 1)$	$\theta \cdot \text{Beta}(\frac{\alpha}{K}, 1)$	$\min(X^{(1)}, \dots, X^{(K)})$
	$F^n$	$F^{n_k}$	$\text{sort}(X^{(1)}, \dots, X^{(K)})$

# The list of distributions we can thin is extensive

Family	Distribution $P_\theta$ , where $X \sim P_\theta$ .	Distribution $Q_\theta^{(k)}$ where $X^{(k)} \stackrel{ind.}{\sim} Q_\theta^{(k)}$ .	Sufficient statistic $T$ (sufficient for $\theta$ )
Natural exponential family (in parameter $\theta$ )	$N(\theta, \sigma^2)$	$N(\epsilon_k \theta, \epsilon_k \sigma^2)$	$\sum_{k=1}^K X^{(k)}$
	Poisson( $\theta$ )	Poisson( $\epsilon_k \theta$ )	
	NegBin( $r, \theta$ )	NegBin( $\epsilon_k r, \theta$ )	
	Binomial( $r, \theta$ )	Binomial( $\epsilon_k r, \theta$ )	
	Gamma( $\alpha, \theta$ )	Gamma( $\epsilon_k \alpha, \theta$ )	
	$N_p(\boldsymbol{\theta}, \Sigma)$	$N_p(\epsilon_k \boldsymbol{\theta}, \epsilon_k \Sigma)$	
	Multinomial $_p(r, \boldsymbol{\theta})$	Multinomial $_p(\epsilon_k r, \boldsymbol{\theta})$	
General exponential family (in parameter $\theta$ )	Gamma( $K/2, \theta$ )	$N(0, \frac{1}{2\theta})$	$\sum_{k=1}^K (X^{(k)})^2$
	Gamma( $K, \theta$ )	Weibull( $\theta^{-\frac{1}{\nu}}, \nu$ )	$\sum_{k=1}^K (X^{(k)})^\nu$
	Beta( $\theta, \beta$ )	Beta( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
	Beta( $\alpha, \theta$ )	Beta( $\frac{1}{K}\alpha, \frac{1}{K}\theta + \frac{k-1}{K}$ )	$(\prod_{k=1}^K (1 - X^{(k)}))^{1/K}$
	Gamma( $\theta, \beta$ )	Gamma( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
	Weibull( $\theta, \nu$ )	Gamma( $\frac{1}{K}, \theta^{-\nu}$ )	$(\sum_{k=1}^K X^{(k)})^{1/\nu}$
	Pareto( $\nu, \theta$ )	Gamma( $\frac{1}{K}, \theta$ )	$\nu \times \text{Exp}(\sum_{k=1}^K X^{(k)})$
	$N(0, \theta)$	Gamma( $\frac{1}{2K}, \frac{1}{2\theta}$ )	$X^2 = \sum_{k=1}^K X^{(k)}$
	$N_K(\theta_1 \mathbf{1}_K, \theta_2 I_K)$	$N(\theta_1, \theta_2)$	sample mean and variance
Truncated support family	Unif( $0, \theta$ )	$\theta \cdot \text{Beta}(\frac{1}{K}, 1)$	$\max(X^{(1)}, \dots, X^{(K)})$
	$\theta \cdot \text{Beta}(\alpha, 1)$	$\theta \cdot \text{Beta}(\frac{\alpha}{K}, 1)$	
	$\theta + \text{Exp}(\lambda)$	$\theta + \text{Exp}(\lambda/K)$	$\min(X^{(1)}, \dots, X^{(K)})$
Non-parametric	$F^n$	$F^{n_k}$	$\text{sort}(X^{(1)}, \dots, X^{(K)})$

# The list of distributions we can thin is extensive

Family	Distribution $P_\theta$ , where $X \sim P_\theta$ .	Distribution $Q_\theta^{(k)}$ where $X^{(k)} \stackrel{ind.}{\sim} Q_\theta^{(k)}$ .	Sufficient statistic $T$ (sufficient for $\theta$ )
Natural exponential family (in parameter $\theta$ )	$N(\theta, \sigma^2)$	$N(\epsilon_k \theta, \epsilon_k \sigma^2)$	
	Poisson( $\theta$ )	Poisson( $\epsilon_k \theta$ )	
	NegBin( $r, \theta$ )	NegBin( $\epsilon_k r, \theta$ )	
	Binomial( $r, \theta$ )	Binomial( $\epsilon_k r, \theta$ )	$\sum_{k=1}^K X^{(k)}$
	Gamma( $\alpha, \theta$ )	Gamma( $\epsilon_k \alpha, \theta$ )	
	$N_p(\boldsymbol{\theta}, \Sigma)$	$N_p(\epsilon_k \boldsymbol{\theta}, \epsilon_k \Sigma)$	
	Multinomial $_p(r, \boldsymbol{\theta})$	Multinomial $_p(\epsilon_k r, \boldsymbol{\theta})$	
General exponential family (in parameter $\theta$ )	Gamma( $K/2, \theta$ )	$N(0, \frac{1}{2\theta})$	$\sum_{k=1}^K (X^{(k)})^2$
	Gamma( $K, \theta$ )	Weibull( $\theta^{-\frac{1}{\nu}}, \nu$ )	$\sum_{k=1}^K (X^{(k)})^\nu$
	Beta( $\theta, \beta$ )	Beta( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
	Beta( $\alpha, \theta$ )	Beta( $\frac{1}{K}\alpha, \frac{1}{K}\theta + \frac{k-1}{K}$ )	$(\prod_{k=1}^K (1 - X^{(k)}))^{1/K}$
	Gamma( $\theta, \beta$ )	Gamma( $\frac{1}{K}\theta + \frac{k-1}{K}, \frac{1}{K}\beta$ )	$(\prod_{k=1}^K X^{(k)})^{1/K}$
	Weibull( $\theta, \nu$ )	Gamma( $\frac{1}{K}, \theta^{-\nu}$ )	$(\sum_{k=1}^K X^{(k)})^{1/\nu}$
	Pareto( $\nu, \theta$ )	Gamma( $\frac{1}{K}, \theta$ )	$\nu \times \text{Exp}(\sum_{k=1}^K X^{(k)})$
Truncated support family	$N(0, \theta)$	Gamma( $\frac{1}{2K}, \frac{1}{2\theta}$ )	$X^2 = \sum_{k=1}^K X^{(k)}$
	$N_K(\theta_1 \mathbf{1}_K, \theta_2 I_K)$	$N(\theta_1, \theta_2)$	sample mean and variance
	$\theta \cdot \text{Beta}(\frac{1}{K}, 1)$	$\theta \cdot \text{Beta}(\frac{\alpha}{K}, 1)$	$\max(X^{(1)}, \dots, X^{(K)})$
Non-parametric	$\theta + \text{Exp}(\lambda)$	$\theta + \text{Exp}(\lambda/K)$	$\min(X^{(1)}, \dots, X^{(K)})$
	$F^n$	$F^{n_k}$	$\text{sort}(X^{(1)}, \dots, X^{(K)})$

We are working on additional extensions and applications

---

The screenshot shows a red header bar with the arXiv logo and navigation links. Below it, a grey header bar indicates the category 'Statistics > Methodology'. The main content area features the title 'Generalized Data Thinning Using Sufficient Statistics' in large bold black font, followed by the authors' names in blue: 'Ameer Dharamshi, Anna Neufeld, Keshav Motwani, Lucy L. Gao, Daniela Witten, Jacob Bien'. The page is identified by the identifier 'arXiv:2303.12931'.

arXiv > stat > arXiv:2303.12931

Search...  
Help | Advanced

Statistics > Methodology

[Submitted on 22 Mar 2023]

**Generalized Data Thinning Using Sufficient Statistics**

Ameer Dharamshi, Anna Neufeld, Keshav Motwani, Lucy L. Gao, Daniela Witten,  
Jacob Bien

arXiv:2303.12931

# Acknowledgements

---



Daniela Witten  
University of Washington



Lucy Gao  
University of British Columbia



Ameer Dharamshi  
University of Washington



Keshav Motwani  
University of Washington



Alexis Battle  
Johns Hopkins



Joshua Popp  
Johns Hopkins



Jacob Bien  
USC

# Questions?

---