



Tuesday, August 6th, 2024

Flexible Synthetic Data Generation with library(tidysynthesis)

Aaron R. Williams

Summary

1. We created a better, more flexible synthesizer for generating FCS synthetic data.

Summary

1. We created a better, more flexible synthesizer for generating FCS synthetic data.
2. Synthetic data are only as good as the synthesizer.

Summary

1. We created a better, more flexible synthesizer for generating FCS synthetic data.
2. Synthetic data are only as good as the synthesizer.
3. Existing tools are limited and difficult to extend.

Summary

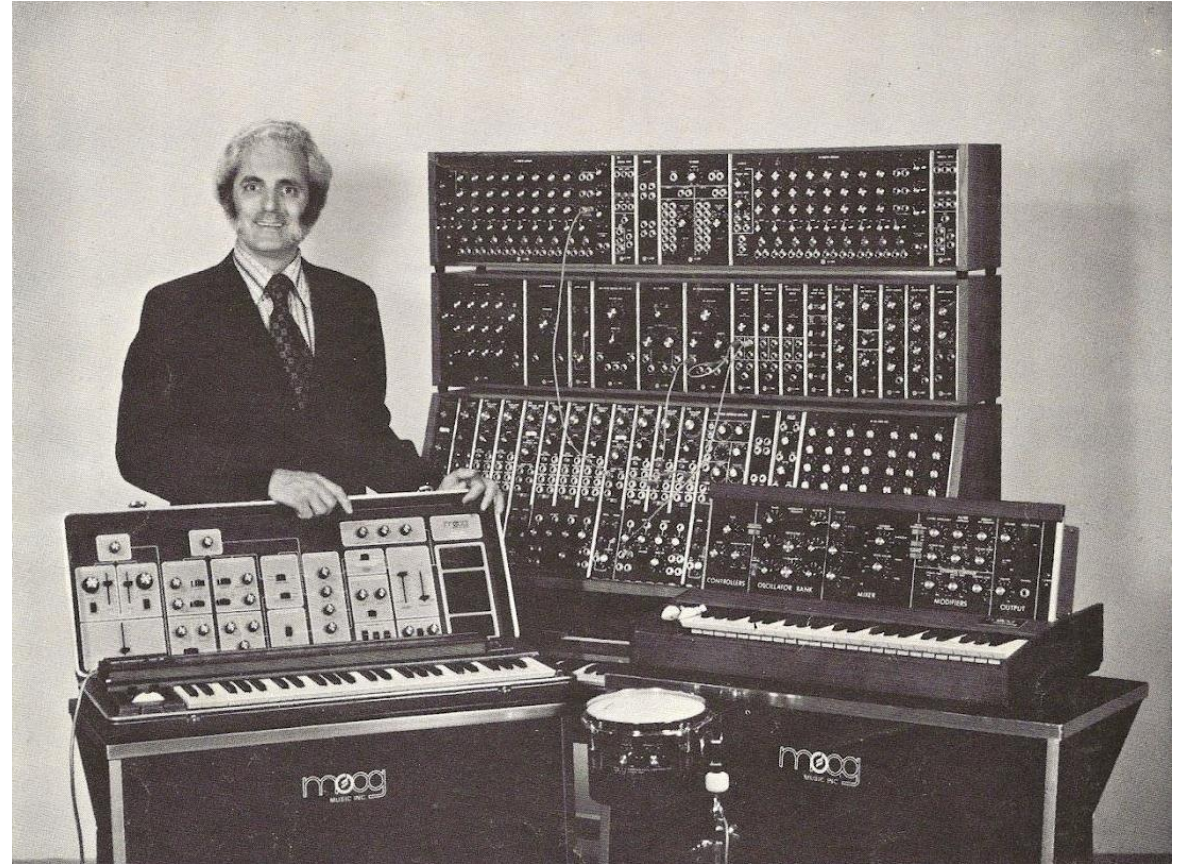
1. We created a better, more flexible synthesizer for generating FCS synthetic data.
2. Synthetic data are only as good as the synthesizer.
3. Existing tools are limited and difficult to extend.
4. `library(tidysynthesis)` is modular, extensible, and builds on `library(tidymodels)`.

Summary

1. We created a better, more flexible synthesizer for generating FCS synthetic data.
2. Synthetic data are only as good as the synthesizer.
3. Existing tools are limited and difficult to extend.
4. `library(tidysynthesis)` is modular, extensible, and builds on `library(tidymodels)`.
5. This works well for generating select synthetic data.

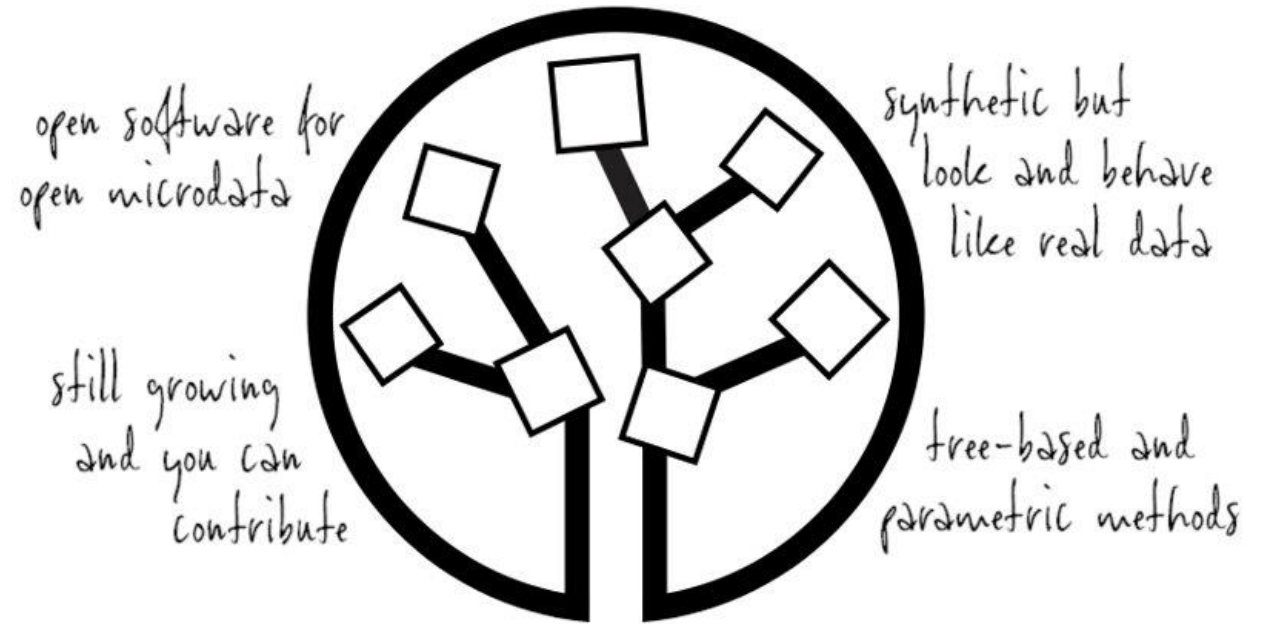
Build a Synthesizer

- Protects the confidentiality of individual information.
- Uses a fully conditional specification to model the joint multivariate distribution.



Existing Functionality is Groundbreaking, but Limited

- Limited to a small set of methods
- Difficult to extend



synthpop

R package for generating synthetic versions of sensitive microdata for statistical disclosure control

R Packages

- **tidysynthesis:** Flexible tools for generating fully and partially synthetic data.
- **syntheval:** Utility and disclosure risk evaluation of synthetic data.

Our Approach

1. Embrace design philosophy from tidyverse and tidymodels
2. Flexible
3. Modular
4. Extensible



library(tidymodels)

- *All the power of library(tidymodels) for data synthesis, concisely, with a few special tools.*
- Full predictive modeling toolkit



Flexibility

- Express different recipes, predictive models, and samplers for each variable
- Hyperparameter tuning
- Additional noise
- Mid-synthesis constraints
- Synthesize missing data
- Parallel computation with futures/furrr

Modular

- Everything is handled through objects with classes
 - Interchangeable objects that act like building blocks
- Robust testing suite
- Manage computation
 - Lazy evaluation and checks that catch errors before computation

Extensibility

- Ability for someone else to add the thing we haven't thought of

Demonstration

Synthetic data

Confidential data

select	species	island	sex	bill_length_mm	...
TRUE	Adelie	Torgersen	male	39.1	...
FALSE	Adelie	Torgersen	female	39.5	...

Synthetic data

select	species	island	sex	bill_length_mm	...
TRUE					
TRUE					

 Synthetic data! 

Synthetic data

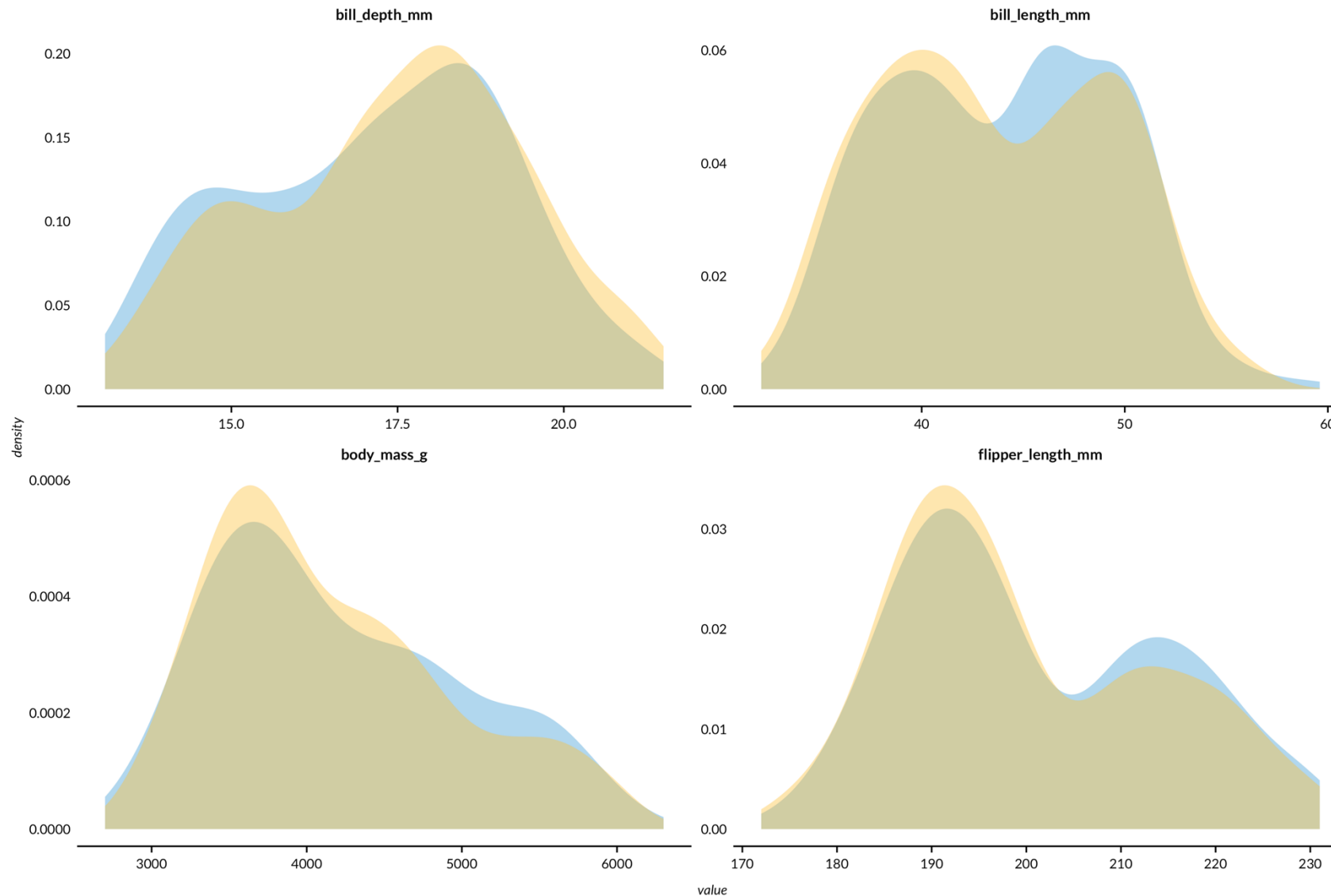
Confidential data

select	species	island	sex	bill_length_mm	...
TRUE	Adelie	Torgersen	male	39.1	...
FALSE	Adelie	Torgersen	female	39.5	...

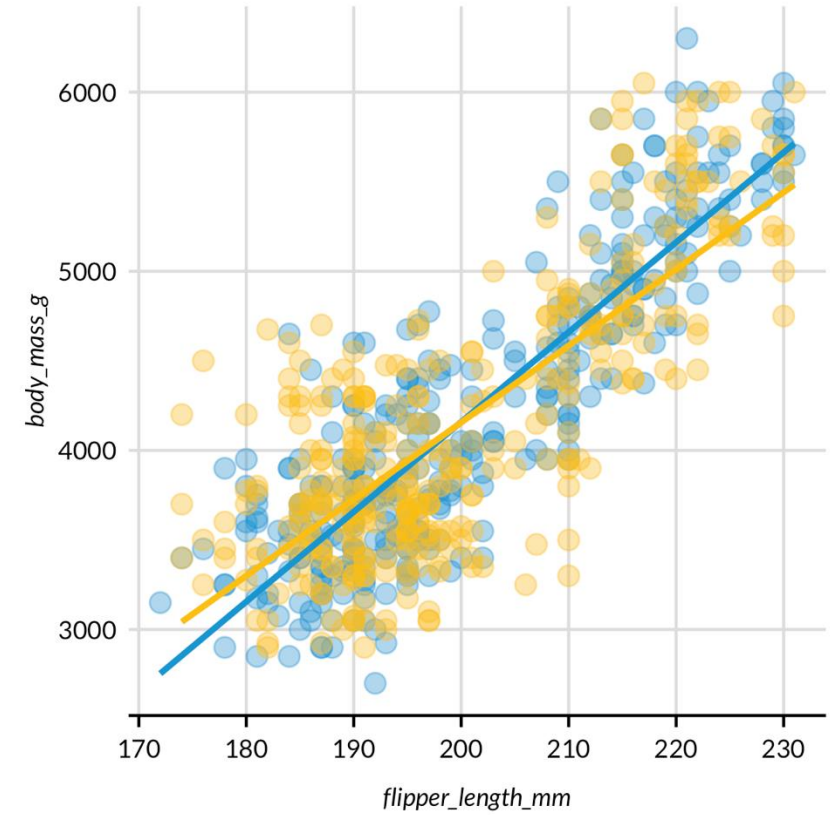
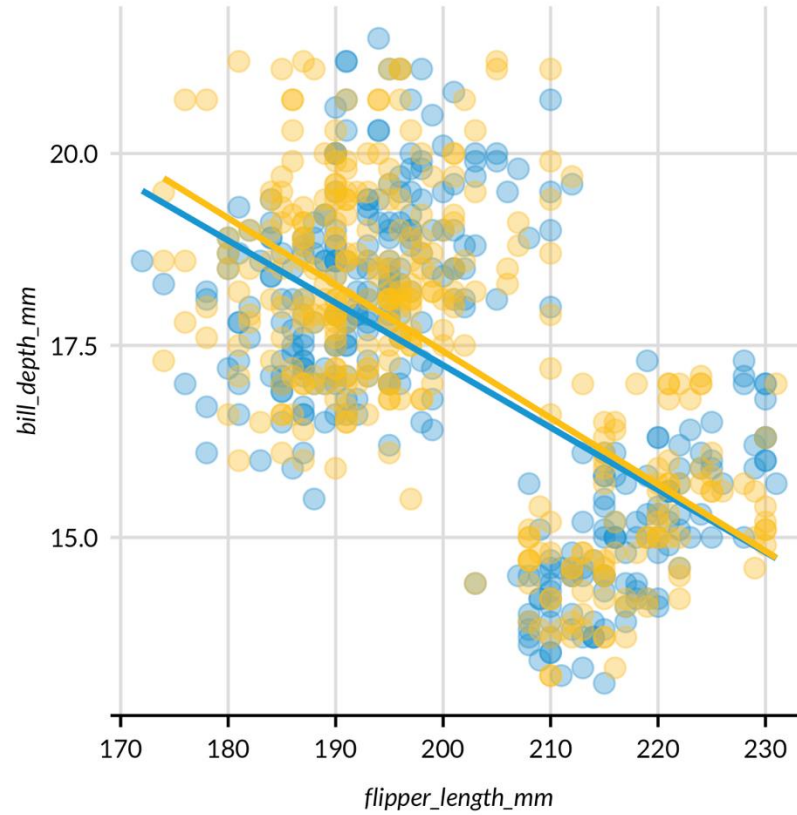
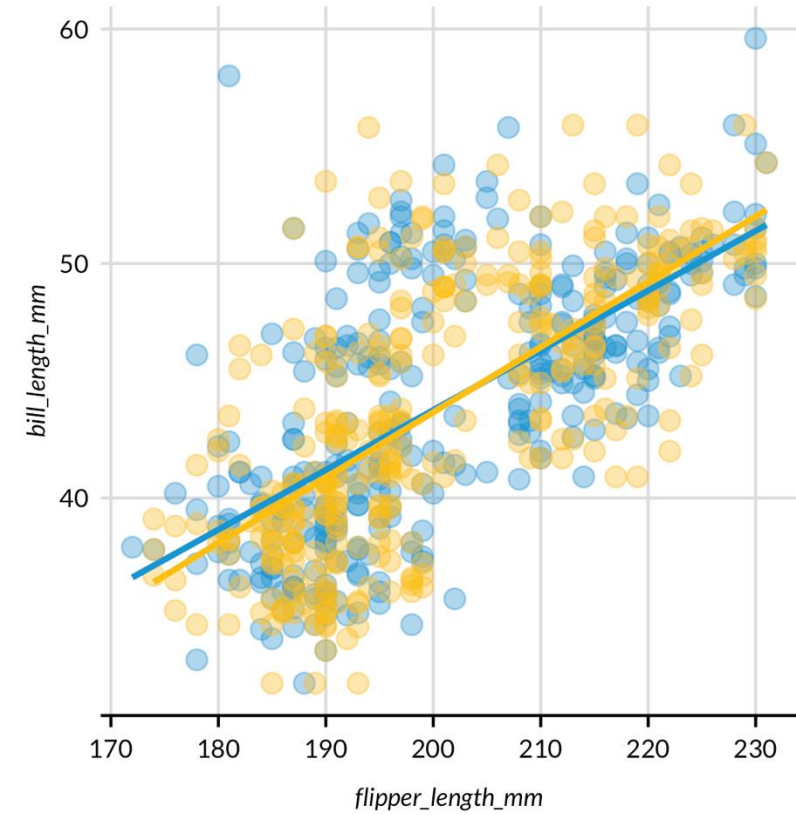
Synthetic data

select	species	island	sex	bill_length_mm	...
TRUE	Chinstrap	Dream	male	48.4	...
TRUE	Gentoo	Biscoe	male	51.4	...

The Synthetic Data are Similar to the Confidential Data



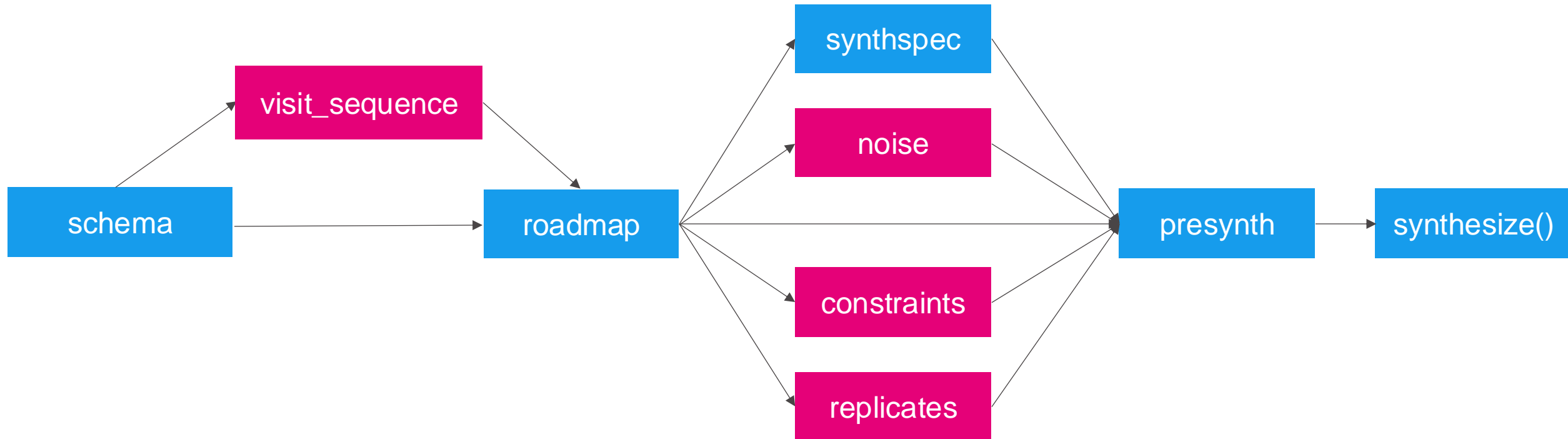
The Synthetic Data are Similar to the Confidential Data



Summary

1. We created a better, more flexible synthesizer for generating FCS synthetic data.
2. Synthetic data are only as good as the synthesizer.
3. Existing tools are limited and difficult to extend.
4. `library(tidysynthesis)` is modular, extensible, and builds on `library(tidymodels)`.
5. This works well for generating select synthetic data.

Workflow



Estimating the multivariate distribution of the data

- Goal is to approximate the empirical multivariate distribution function for the data
- Joint multivariate probability distribution can be represented as the product of sequential, conditional probability distributions:

$$f(Y_1, Y_2, \dots, Y_k | \theta_1, \theta_2, \dots, \theta_k) =$$

$$f_1(Y_1 | \theta_1) \cdot f_2(Y_2 | Y_1, \theta_2) \cdots f_k(Y_k | Y_1, Y_2, \dots, Y_{k-1}, \theta_k)$$

- where Y_i the variables and θ_i are vectors of model parameters