

# Supplementary materials to TITLE

Martin Papenberg

## Motivating Example

The synthetic data set that we generated to resemble our actual application consisted of 370 samples from 191 unique patients. Samples belonging to the same patient were required to be assigned to the same batch. Most patients ( $n = 100$ ) were included with 1 sample. The remaining 270 patients provided more than one sample: 47 patients provided 2 samples, 13 patients provided 3 samples, 24 patients provided 4 samples, 3 patients provided 5 samples, 3 patients provided 6 samples, and 1 patient provided 8 samples. We required the samples to be assigned to 16 equal-sized batches. As a perfect split was not possible using this constellation, we created 2 batches containing 24 samples and 14 batches containing 23 samples. We strived for balance among the batches with regard to two numeric variables (age, BMI) and several categorical variables (such as ethnicity, disease stage, and cycle phase). Table 1 illustrates the results of three batch assignments using three of the variables (age, BMI, ethnicity) for illustrative purposes. The anticlustering assignments were implemented using the `anticlust` R package, and the code and data to reproduce the assignment can be retrieved via XXX (**TODO**). The first assignment is based on “standard” anticlustering, which ignores the must-link constraints and assumes that all samples can be assigned independently, with the objective of maximum similarity among batches. This unrestricted anticlustering led to the most pronounced balance among the 16 batches: Average age ranged between 30 and 30.17 years, average BMI ranged between 29.38 and 29.43, and the percentage of samples belonging to hispanic patients ranged between 50 and 54%. The second assignment implemented the must-link restrictions, ensuring that samples belonging to the same patient were assigned to the same batch, while still optimizing balance via anticlustering. The constrained assignment reduced the overall similarity among batches slightly, but the variables were arguably still rather well balanced: Among the 16 batches, average age ranged between 29.91 and 30.35 years, average BMI ranged between 29.3 and 29.57, and the percentage of samples belonging to hispanic patients ranged between 50 and 57%. The third assignment used completely random allocation of samples to batches. Random assignment is often considered as the intuitive and simple method for allocation when no other method is available. However, it does not necessarily lead to similarity among batches (Yan et al. 2012; Papenberg and Klau 2021). In this application, random assignment by far led to the worst balance: Average age ranged between 28.04 and 31.96 years, average BMI ranged between 28.04 and 31.43, and the percentage of samples belonging to hispanic patients ranged between 30 and 65%.

In the next section, we will explain the methodology of anticlustering that led to the batch assignments illustrated in Table 1. We will reiterate the general anticlustering methodology first, before outlining how we included must-link constraints with anticlustering. The implementation of must-link constraints with anticlustering is a novel contribution of the current paper.

## Methods

Anticlustering is an optimization method that is characterized by (a) an objective function that quantifies the balance among batches, and (b) an algorithm that conducts the batch assignment in such a way that balance among batches is maximized. Anticlustering owes its name to the fact that the objective functions it uses are the reversal of criteria used in cluster analysis. For example, Späth (1986) already recognized that by maximizing instead of minimizing the k-means criterion (the “variance”), he was able to create groups that are similar to each other, and presented it as an improvement over the more intuitive random assignment

Table 1: Balance among 16 batches after implementing three different assignment procedures

Batch	Unconstrained Anticlustering			Constrained Anticlustering			Random Assignment		
	Mean Age	Mean BMI	% Hispanic	Mean Age	Mean BMI	% Hispanic	Mean Age	Mean BMI	% Hispanic
1	30.00	29.42	54	30.17	29.38	50	30.62	30.17	42
2	30.08	29.38	50	30.12	29.38	50	28.04	29.29	62
3	30.17	29.43	52	30.13	29.39	52	29.43	29.04	52
4	30.13	29.39	52	30.09	29.39	52	30.17	29.00	57
5	30.13	29.43	52	30.00	29.43	52	31.78	29.78	52
6	30.13	29.39	52	30.35	29.52	52	30.00	29.61	39
7	30.13	29.39	52	30.09	29.43	52	28.78	30.04	61
8	30.13	29.43	52	30.17	29.30	52	29.43	29.61	57
9	30.17	29.43	52	30.00	29.52	52	29.74	28.13	65
10	30.13	29.43	52	30.04	29.43	52	31.65	28.26	35
11	30.13	29.43	52	30.30	29.30	52	31.52	30.48	57
12	30.17	29.43	52	30.22	29.57	57	29.83	28.39	30
13	30.13	29.39	52	30.35	29.43	52	28.87	28.04	57
14	30.13	29.43	52	30.09	29.43	52	28.52	30.13	65
15	30.17	29.43	52	29.91	29.39	52	31.83	29.26	52
16	30.17	29.43	52	30.09	29.39	52	31.96	31.43	52

(Steinley, 2006). Brusco et al. (2020) recognized that other objective functions known from cluster analysis can also be implemented in the context of anticlustering.

In our application, we optimized the diversity objective to maximize similarity among batches. While the diversity is technically a measure of within-batch heterogeneity, its maximization simultaneously leads to minimal difference between the distribution of the input variables among batches (Feo & Khellaf 1990; cf. Papenberg, 2024). Papenberg and Klau (2021) referred to the maximization of the diversity as anticluster editing because the minimization of the diversity is also well-known from the area of cluster analysis—under the term “cluster editing” (Shamir et al., 2004; Böcker et al., 2011). The diversity is computed on the basis of a measure of pairwise dissimilarity among samples. In particular, it is defined as the overall sum of all dissimilarities among samples that are assigned to the same batch (Brusco et al., 2020). Hence, the diversity is not directly computed on the basis of the samples’ features, but instead it relies on distance measure that is computed on the basis of the features, for each pair of samples. In the context of anticlustering, the Euclidean distance is the most common measure that translates features to pairwise dissimilarities (Gallego et al. 2013; Papenberg & Klau, 2021). However, using other distance measures such as the squared Euclidean distance is also possible (Brusco et al., 2020). The Euclidean distance is defined as

$$d(x, y) = \sqrt{\sum_{i=1}^M (x_i - y_i)^2}$$

When samples are described by two features, the Euclidean distance corresponds to the geometric, “straight line” distance between two points in a two-dimensional space; more similar data points are closer to each other (see Figure 1). For categorical variables, we use binary coding before including them in the computation of the Euclidean distance (see Table 2). An anticlustering algorithm assigns samples to batches in such a way that the objective function—here, the diversity—is maximized. Anticlustering usually employs heuristic optimization algorithms (Yang et al., 2022). While heuristics generally provide satisfying results in the context of anticlustering (Papenberg & Klau, 2021), they do not guarantee to find the globally best assignment among all possibilities. In principle, enumerating all possible assignments is a valid strategy to obtain an optimal assignment. However, this approach quickly becomes impossible due to an exponential growth of the way in which assignments can be conducted (Papenberg & Klau, 2021). Moreover, because anticlustering problems (with the exception of some special cases) are also NP-hard (Feo & Khellaf, 1990), there is probably no algorithm that identifies the globally best assignment without considering all possibilities (at least in the worst case). In practice, heuristics are therefore indispensable. In the following, we present an anticlustering

heuristic that is generally applicable to large data sets, and outline how we included must-link constraints with this algorithm.

We used an heuristic exchange algorithm to maximize the diversity (Späth 1986; Papenberg & Klau, 2021; Weitz & Lakshminarayanan, 1998). It consists of two steps: an initialization step and an optimization step. As initialization, it randomly assigns samples to equal-sized batches. In principle, unequal-sized batches would also be possible, but equal-sized batches were required in the current application (and in general, this requirement is most common. After initialization, the algorithm selects the first sample and checks how the diversity would change if the sample were swapped with each sample that is currently assigned to a different batch. After simulating each exchange, it realizes the one exchange that increases the diversity the most. It does not conduct an exchange if no improvement in diversity is possible. This procedure is repeated for each sample and it terminates after the last sample was processed. The procedure might also restart at the first element and reiterate through all samples until no exchange leads to an improvement any more, i.e., until a local maximum is found. In anticlust, we also implemented this local maximum search, which corresponds to the algorithm LCW by Weitz and Lakshminarayanan (1998).<sup>7</sup> For better results, it is also possible to restart the search algorithm multiple times using different (random) initializations (Späth 1986). For the anticlustering assignment illustrated in Table 1, we employed 10 repetitions of the local maximum search LCW, using the squared Euclidean distance as measure of pairwise dissimilarity. To ensure comparable weight of the three features, we also applied a standardization of the input variables before the distance was computed. Using a standardization is recommended if the variables differ strongly in their ranges (Papenberg, 2024).

Table 2: Illustrates the recoding of the categorical variable Race using four binary variables.

Race	Black	Native American	Pacific Islander	White
Black	1	0	0	0
Pacific Islander	0	0	1	0
Native American	0	1	0	0
White	0	0	0	1
Asian	0	0	0	0

### Including Must-Link Constraints with Anticlustering

In our application, samples belonging to the same patient were required to be assigned to the same batch. We refer to a set of samples that must be assigned to the same batch as a must-link clique. In our application, not all samples were part of a clique. For the synthetic data set that resembled our actual application, we simulated 370 samples from 191 unique patients. To include the must-link constraints with anticlustering, we basically use a downscaled data set where each unique patient—but not each single sample—constitutes a unit in the anticlustering process. Hence, in our application, the effective sample size was 191 instead of 370. Some adjustments of the exchange method are required to ensure (a) we still obtain a valid partitioning regarding the size constraints (equal-sized batches) and (b) the diversity is computed correctly during optimization. We therefore had to adjust both the initialization phase as well as the optimization phase of the exchange algorithm.

During initialization, we first assign all samples to a batch that are part of a clique (270 samples). Each clique must be assigned completely to one of the batches and samples within a clique must not be split apart. At the same time, the maximum capacity of each batch must not be exceeded. Using this conceptualization, the initialization step corresponds to a bin packing problem, which is one of the classical NP complete problems in computer science (Garey & Johnson, 1979). That is, we assign a weight to each clique, corresponding to the number of samples it contains. When filling batches, the sum of the weights of the cliques in each batch must not exceed its capacity. Many optimal and heuristic algorithms have been developed to address such a bin packing problem. As the default method, we use a randomized first fit heuristic to fill batches: For each must-link clique, we iterate through all batches in random order and assign it to the first batch where it fits.

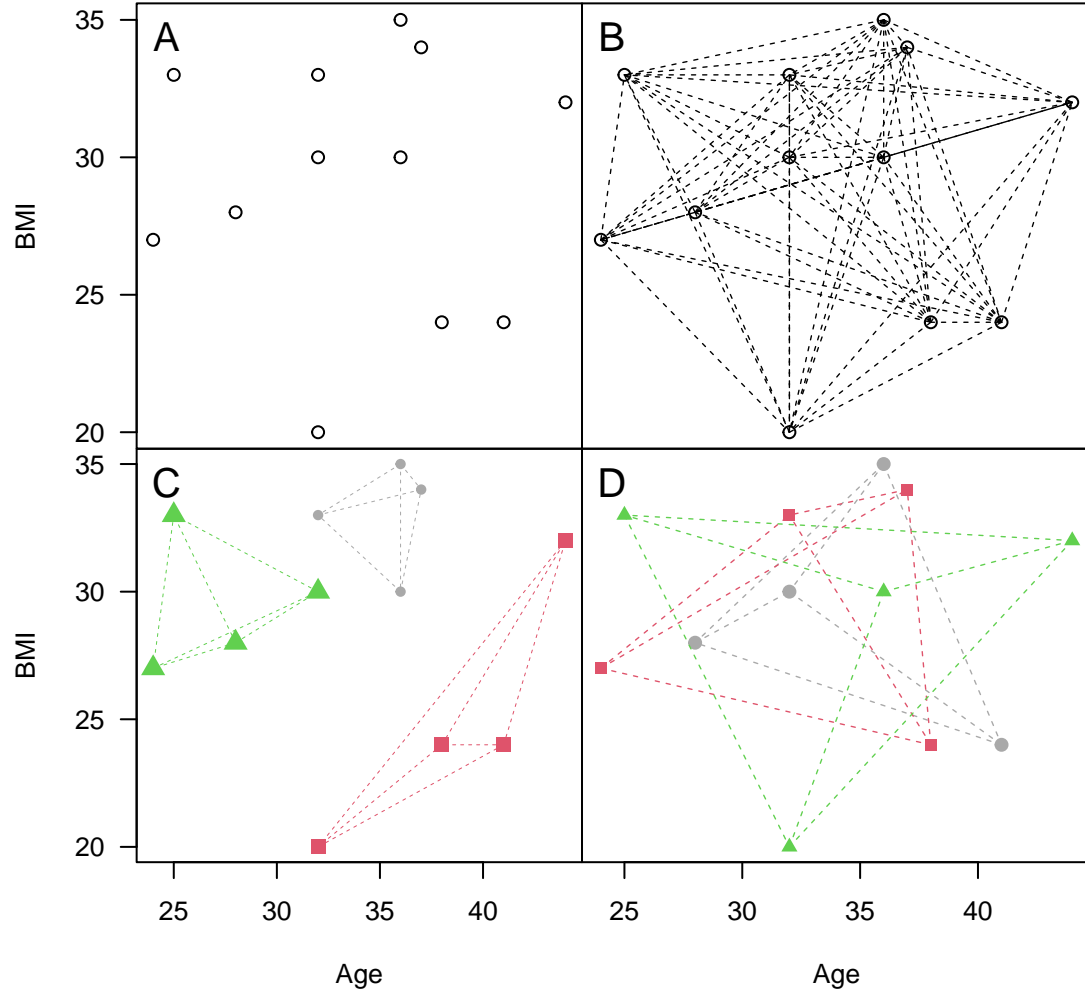


Figure 1: Illustrates the conversion from numeric features to Euclidean distance, and (anti)clustering assignments based on minimum and maximum diversity using the Euclidean distance. Panel A illustrates the BMI and age of twelve women in our synthetic data in a scatter plot. Panel B represents the Euclidean distances between features as a straight line in the two-dimensional space. The Euclidean distance is proportional to the length of the connecting lines in panel B. Panel C illustrates a clustering assignment of the 12 data points to  $K = 3$  equal-sized groups via *minimum* diversity. Panel D illustrates an anticlustering assignment of the 12 data points to  $K = 3$  equal-sized groups via *maximum* diversity. The diversity is computed as the sum of within-(anti)cluster distances, which are highlighted in Panel C and Panel D through connecting lines. Maximizing the diversity simultaneously leads to similar distribution of the input features among batches.

The process is expected to evenly distribute the must-link cliques among batches. This random component is particularly useful if we use multiple restarts of the optimization algorithm. After assigning the must-link cliques to batches, the remaining samples can be assigned randomly to fill the remaining space. Note that our randomized first fit algorithm is a heuristic that may not find an assignment of must-link groups to batches even if one is theoretically available. If the heuristic indicates that the batches cannot hold the must-link groups, we therefore use an optimal algorithm based on integer linear programming as a fallback option, which allows us to verify if the constraints really cannot be fulfilled. To this end, we implemented an adaptation of the standard bin packing ILP model by Martello and Toth, (1990, 221). It is given as

$$\text{minimize} \quad \sum_{1 \leq i \leq K} \sum_{1 \leq j \leq n} x_{ij} \quad (1)$$

$$\text{subject to} \quad \sum_{j=1}^n w_j x_{ij} \leq c_i \quad i = (1, \dots, K) \quad (2)$$

$$\sum_{i=1}^K x_{ij} = 1 \quad j = (1, \dots, n) \quad (3)$$

$$x_{ij} \in \{0, 1\} \quad i = (1, \dots, K), j = (1, \dots, n) \quad (4)$$

The number of must-link cliques is given by  $n$ . The model has decision variables  $x_{ij}$  to encode whether clique  $j$  ( $j = 1, \dots, n$ ) is assigned to batch  $i$  ( $i = 1, \dots, K$ ). It uses  $K$  values  $c_i$  to represent the capacity of each batch; in the default case of equal-sized batches, we have  $c_i = \frac{N}{K}$  for each batch. It uses  $n$  values  $w_j$  to encode the weight of each clique, i.e., the number of samples it represents that must be assigned to the same batch in order to fulfil the must-link constraints. Constraint (2) realizes that the weight of each batch is not exceeded; constraint (3) realizes that each clique is assigned to exactly one batch. Note that during the initialization step that assigns cliques to batches, we only need to test if the constraints (2) and (3) can be fulfilled at all; any feasible assignment is equally valid. For this reason, the objective function (1) is chosen to be constant for each assignment that satisfies the constraints ( $n$ ). It does not actually contribute to solving the problem, and the model only test if the must-link constraints can be fulfilled.

Like the initialization step, the optimization step of the exchange algorithm uses a downscaled data set where each patient rather than each sample corresponds to a unit of analysis. However, to obtain valid results according to the original data set (which incorporates all samples), an adjustment to the data input is needed. In particular, we must change the matrix of pairwise dissimilarities: To obtain a reduced distance matrix that preserves all information of the original distance matrix, we sum up all pairwise distances between samples in different cliques. In the context of maximizing the diversity, this transformation sufficiently preserves the relevant information in the original distance matrix (Böcker et al., 2011). Using the initial assignment and the reduced distance matrix, we apply the same exchange algorithm as we described for the unrestricted anticlustering application, with one adjustment: During the exchange process, we only exchange cliques of the same size (e.g., patients providing the same number of samples) to ensure that the cardinality constraints are respected throughout (i.e., usually equal-sized batches).

### Optimal anticlustering using must-link constraints

As mentioned above due the computational complexity of anticlustering, batch assignment problems are usually tackled using heuristic algorithms. Still, for some problem constellations, in particular when  $N$  is not large, it is possible to employ optimal algorithms that find the globally best batch assignment. Papenberg and Klau (2021) presented an ILP model to find globally optimal batch assignments for the diversity objective. It can be used to solve problem instances of up to about  $N = 30$  in an acceptable running time; Schulz (2022; ). In this paper, we extend the model by Papenberg and Klau (2021) by allowing it to include must-link constraints. The extension is actually quite straight forward: To induce must-link constraints in the context of an optimal algorithm, it is sufficient to adjust the distance matrix used as input. Whenever the pairwise distance between two samples is set to  $\infty$ —and the set of must-link constraints can be fulfilled—any globally optimal assignment will place these samples in the same batch, because the objective value associated with such an assignment is necessarily better than that of an assignment that places them in different batches. This

feature of adjusting the data input has long been used in the context of optimal algorithms for cluster editing (Böcker et al. 2011). Including must-link constraints in the anticlustering problem specification increases the number of samples that can be processed. In our online supplementary materials, we show that up to about 40-50 samples can be processed in a time limit of 1800 seconds on a personal computer; after that, the running time is expected to increase exponentially with increasing number of samples.

## Simulation Study

For the simulation, we generated 2000 data sets. Data sets were processed via (a) OSAT, (b) PSBA, (c) unconstrained anticlustering and (d) anticlustering subject to must-link constraints. Because the OSAT method is only applicable to categorical variables, we used categorical variables in our simulation. For anticlustering and for PSBA, the categorical variables were binary coded before the methods were applied (see Table 2). For each data set, we randomly determined the number of categorical variables (2-5), the number of classes per variable (2-5; the distribution of classes was uniform), the total sample size  $N$  (between 50 and 500) and the number of batches  $K$  (2, 4, or 10 equal-sized batches). PSBA was only applied for  $K = 2$  and  $K = 4$  ( $n = 1375$  data sets) because the authors' implementation only allows the assignment to a maximum of four batches. Must-link constraints were generating a random integer between 1 and  $N$  (with equal probability), and using the resulting numbers as must-link grouping variables. This rule resulted in a distribution of constraints that resembled our motivating application: 58% of all elements had no must-link partner; 29% had 1 must-link partner; 10% had 2 must-link partners, and 3% had 3 or more must-link partners.

For each simulation run, we computed  $\chi^2$ -tests to assess the imbalance among batches for each of the 2-5 variables, for each of the three competing methods. We stored the  $p$ -value associated with each test. A higher  $p$ -value indicates that there is less imbalance among batches, i.e., that the batches are more similar. The simulation revealed that in 83% of all variable comparisons, balance among batches was better when using anticlustering as compared to OSAT. Balance was equal in 16% of all comparisons, and only in 1% OSAT outperformed **anticlust**. In 51% of all variable comparisons, balance was better when using anticlustering as compared to PSBA. Balance was equal in 16% of all comparisons, and in 34% PSBA outperformed anticlustering. Figure 2 illustrates the average  $p$ -values in dependence of the number of variables ( $M$ ) and the number of batches ( $K$ ). The online supplement also includes additional Figures illustrating average  $p$ -values in dependence of the other factors that varied in the simulation (the sample size  $N$  and the number of classes per categorical variables; **TODO**). Figure 2 shows that increasing the number of variables posed severe challenges for OSAT, but hardly affected anticlustering. PSBA also showed decreased performance with an increasing number of variables, but less so than OSAT.

The anticlustering assignment that was subjected to must-link constraints on average achieved 99.8% of the objective value of the unconstrained assignment. Hence, must-link constraints are not only desirable from a user's point of view, but they also do not decrease batch balance considerably; in 57% of all cases, balance was not at all reduced by the constraints. Remarkably, the constrained anticlustering assignment led to better balance than the OSAT and PSBA assignments that did not employ any constraints (see Figure 2).

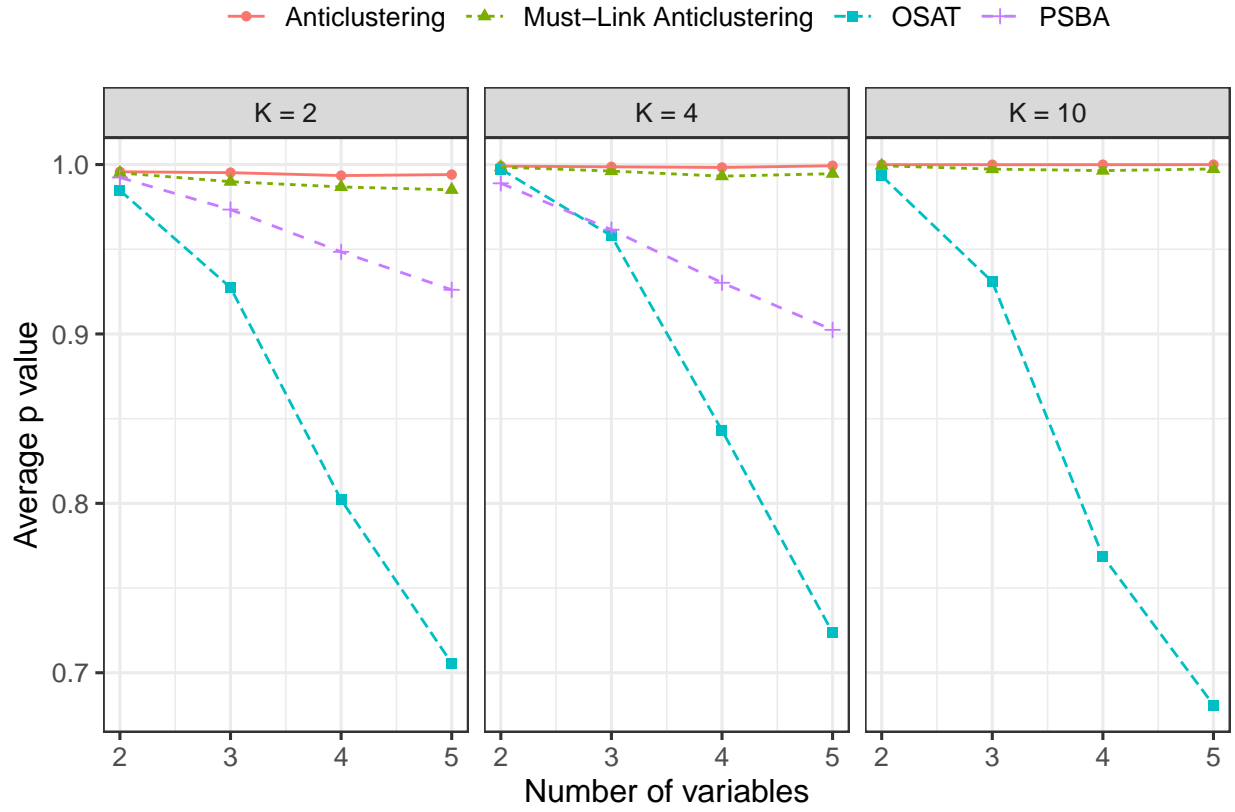


Figure 2: Average  $p$  values in dependence of the number of batches and the number of variables. Higher  $p$  values indicate better balance. Anticlustering maintained a comparable level of balance in all conditions. OSAT's performance decreased with increasing number of variables most strongly.

## References

- Martello, S, and P Toth. 1990. *Knapsack Problems: Algorithms and Computer Implementations*. Wiley.
- Papenberg, Martin, and Gunnar W Klau. 2021. “Using Anticlustering to Partition Data Sets into Equivalent Parts.” *Psychological Methods* 26 (2): 161–74. <https://doi.org/10.1037/met0000301>.
- Yan, Li, Changxing Ma, Dan Wang, Qiang Hu, Maochun Qin, Jeffrey M Conroy, Lara E Sucheston, et al. 2012. “OSAT: A Tool for Sample-to-Batch Allocations in Genomics Experiments.” *BMC Genomics* 13: 1–7.