
Cluster Stability for Finite Samples

Ohad Shamir[†] and Naftali Tishby^{†‡}

[†] School of Computer Science and Engineering

[‡] Interdisciplinary Center for Neural Computation

The Hebrew University

Jerusalem 91904, Israel

{ohadsh,tishby}@cs.huji.ac.il

Abstract

Over the past few years, the notion of stability in data clustering has received growing attention as a cluster validation criterion in a sample-based framework. However, recent work has shown that as the sample size increases, any clustering model will usually become asymptotically stable. This led to the conclusion that stability is lacking as a theoretical and practical tool. The discrepancy between this conclusion and the success of stability in practice has remained an open question, which we attempt to address. Our theoretical approach is that stability, as used by cluster validation algorithms, is similar in certain respects to measures of generalization in a model-selection framework. In such cases, the model chosen governs the convergence rate of generalization bounds. By arguing that these rates are more important than the sample size, we are led to the prediction that stability-based cluster validation algorithms should not degrade with increasing sample size, despite the asymptotic universal stability. This prediction is substantiated by a theoretical analysis as well as some empirical results. We conclude that stability remains a meaningful cluster validation criterion over finite samples.

1 Introduction

Clustering is one of the most common tools of unsupervised data analysis. Despite its widespread use and an immense amount of literature, distressingly little is known about its theoretical foundations [14]. In this paper, we focus on sample based clustering, where it is assumed that the data to be clustered are actually a sample from some underlying distribution.

A major problem in such a setting is assessing cluster validity. In other words, we might wish to know whether the clustering we have found actually corresponds to a meaningful clustering of the underlying distribution, and is not just an artifact of the sampling process. This problem relates to the issue of model selection, such as determining the number of clusters in the data or tuning parameters of the clustering algorithm. In the past few years, cluster stability has received growing attention as a criterion for addressing this problem. Informally, this criterion states that if the clustering algorithm is repeatedly applied over independent samples, resulting in 'similar' clusterings, then these clusterings are statistically significant. Based on this idea, several cluster validity methods have been proposed (see [9] and references therein), and were shown to be relatively successful for various data sets in practice.

However, in recent work, it was proven that under mild conditions, stability is asymptotically fully determined by the behavior of the objective function which the clustering algorithm attempts to optimize. In particular, the existence of a unique optimal solution for some model choice implies stability as sample size increase to infinity. This will happen regardless of the model fit to the data. From this, it was concluded that stability is not a well-suited tool for model selection in clustering. This left open, however, the question of why stability is observed to be useful in practice.

In this paper, we attempt to explain why stability measures should have much wider relevance than what might be concluded from these results. Our underlying approach is to view stability as a measure of generalization, in a learning-theoretic sense. When we have a ‘good’ model, which is stable over independent samples, then inferring its fit to the underlying distribution should be easy. In other words, stability should ‘work’ because stable models generalize better, and models which generalize better should fit the underlying distribution better. We emphasize that this idea in itself is not novel, appearing explicitly and under various guises in many aspects of machine learning. The novelty in this paper lies mainly in the predictions that are drawn from it for clustering stability.

The viewpoint above places emphasis on the nature of stability for *finite* samples. Since generalization is meaningless when the sample is infinite, it should come as no surprise that stability displays similar behavior. On finite samples, the generalization uncertainty is virtually always strictly positive, with different model choices leading to different convergence rates towards zero for increasing sample size. Based on the link between stability and generalization, we predict that on realistic data, all risk-minimizing models asymptotically become stable, but the *rates of convergence* to this ultimate stability differ. In other words, an appropriate scaling of the stability measures will make them independent of the actual sample size used. Using this intuition, we characterize and prove a mild set of conditions, applicable in principle to a wide class of clustering settings, which ensure the relevance of cluster stability for arbitrarily large sample sizes. We then prove that the stability measure used in previous work to show negative asymptotic results on stability, actually allows us to discern the ‘correct’ model, regardless of how large is the sample, for a certain simple setting. Our results are further validated by some experiments on synthetic and real world data.

2 Definitions and notation

We assume that the data sample to be clustered, $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, is produced by sampling instances i.i.d from an underlying distribution \mathcal{D} , supported on a subset \mathcal{X} of \mathbb{R}^n . A *clustering* C_D for some $D \subseteq \mathcal{X}$ is a function from $D \times D$ to $\{0, 1\}$, defining an equivalence relation on D with a finite number of equivalence classes (namely, $C_D(\mathbf{x}_i, \mathbf{x}_j) = 1$ if \mathbf{x}_i and \mathbf{x}_j belong to the same cluster, and 0 otherwise). For a clustering $C_{\mathcal{X}}$ of the instance space, and a finite sample S , let $C_{\mathcal{X}}|_S$ denote the functional restriction of $C_{\mathcal{X}}$ on $S \times S$.

A *clustering algorithm* A is a function from any finite sample $S \subseteq \mathcal{X}$, to some clustering $C_{\mathcal{X}}$ of the instance space¹. We assume the algorithm is driven by optimizing an objective function, and has some user-defined parameters Θ . In particular, A_k denotes the algorithm A with the number of clusters chosen to be k .

Following [2], we define the *stability* of a clustering algorithm A on finite samples of size m as:

$$\text{stab}(A, \mathcal{D}, m) = \mathbb{E}_{S_1, S_2} d_{\mathcal{D}}(A(S_1), A(S_2)), \quad (1)$$

where S_1 and S_2 are samples of size m , drawn i.i.d from \mathcal{D} , and $d_{\mathcal{D}}$ is some ‘dissimilarity’ function between clusterings of \mathcal{X} , to be specified later.

Let ℓ denote a *loss function* from any clustering C_S of a finite set $S \subseteq \mathcal{X}$ to $[0, 1]$. ℓ may or may not correspond to the objective function the clustering algorithm attempts to optimize, and may involve a global quality measure rather than some average over individual instances. For a fixed sample size, we say that ℓ obeys the *bounded differences property* (see [11]), if for any clustering C_S it holds that $|\ell(C_S) - \ell(C_{S'})| \leq a$, where a is a constant, and $C_{S'}$ is obtained from C_S by replacing at most one instance of S by any other instance from \mathcal{X} , and clustering it arbitrarily.

A *hypothesis class* H is defined as some set of clusterings of \mathcal{X} . The *empirical risk* of a clustering $C_{\mathcal{X}} \in H$ on a sample S of size m is $\ell(C_{\mathcal{X}}|_S)$. The *expected risk* of $C_{\mathcal{X}}$, with respect to samples S of size m , will be defined as $\mathbb{E}_S \ell(C_{\mathcal{X}}|_S)$. The problem of generalization is how to estimate the expected risk, based on the empirical data.

¹Many clustering algorithms, such as spectral clustering, do not induce a natural clustering on \mathcal{X} based on a clustering of a sample. In that case, we view the algorithm as a two-stage process, in which the clustering of the sample is extended to \mathcal{X} through some uniform extension operator (such as assigning instances to the ‘nearest’ cluster in some appropriate sense).

3 A Bayesian framework for relating stability and generalization

The relationship between generalization and various notions of stability is long known, but has been dealt with mostly in a supervised learning setting (see [3][5] [8] and references therein). In the context of unsupervised data clustering, several papers have explored the relevance of statistical stability and generalization, separately and together (such as [1][4][14][12]). However, there are not many theoretical results quantitatively characterizing the relationship between the two in this setting. The aim of this section is to informally motivate our approach, of viewing stability and generalization in clustering as closely related.

Relating the two is very natural in a Bayesian setting, where clustering stability implies an ‘unsurprising’ posterior given a prior, which is based on clustering another sample. Under this paradigm, we might consider ‘soft clustering’ algorithms which return a distribution over a measurable hypothesis class H , rather than a specific clustering. This distribution typically reflects the likelihood of a clustering hypothesis, given the data and prior assumptions. Extending our notation, we have that for any sample S , $A(S)$ is now a distribution over H . The empirical risk of such a distribution, with respect to sample S' , is defined as $\ell(A(S)|_{S'}) = \mathbb{E}_{C_{\mathcal{X}} \sim A(S)} \ell(C_{\mathcal{X}}|_{S'})$.

In this setting, consider for example the following simple procedure to derive a clustering hypothesis distribution, as well as a generalization bound: Given a sample of size $2m$ drawn i.i.d from \mathcal{D} , we randomly split it into two samples S_1, S_2 each of size m , and use A to cluster each of them separately. Then we have the following:

Theorem 1. *For the procedure defined above, assume ℓ obeys the bounded differences property with parameter $1/m$. Define the clustering distance $d_{\mathcal{D}}(\mathcal{P}, \mathcal{Q})$ in Eq. (1), between two distributions \mathcal{P}, \mathcal{Q} over the hypothesis class H , as the Kullback-Leibler divergence $D_{KL}[\mathcal{Q}||\mathcal{P}]^2$. Then for a fixed confidence parameter $\delta \in (0, 1)$, it holds with probability at least $1 - \delta$ over the draw of samples S_1 and S_2 of size m , that*

$$\mathbb{E}_S \ell(A(S_2)|_{S_1}) - \ell(A(S_2)|_{S_2}) \leq \sqrt{\frac{d_{\mathcal{D}}(A(S_1), A(S_2)) + \ln(m/\delta) + 2}{2m - 1}}.$$

The theorem is a straightforward variant of the PAC-Bayesian theorem [10]. Since the loss function is not necessarily an empirical average, we need to utilize McDiarmid’s bound for random variables with bounded differences, instead of Hoeffding’s bound. Other than that, the proof is identical, and is therefore omitted.

This theorem implies that the more stable is the Bayesian algorithm, the tighter the expected generalization bounds we can achieve. In fact, the ‘expected’ magnitude of the high-probability bound we will get (over drawing S_1 and S_2 and performing the procedure described above) is:

$$\begin{aligned} \mathbb{E}_{S_1, S_2} \sqrt{\frac{d_{\mathcal{D}}(A(S_1), A(S_2)) + \ln(m/\delta) + 2}{2m - 1}} &\leq \sqrt{\frac{\mathbb{E}_{S_1, S_2} d_{\mathcal{D}}(A(S_1), A(S_2)) + \ln(m/\delta) + 2}{2m - 1}} \\ &= \sqrt{\frac{\text{stab}(A, \mathcal{D}, m) + \ln(m/\delta) + 2}{2m - 1}}. \end{aligned}$$

Note that the only model-dependent quantity in the expression above is $\text{stab}(A, \mathcal{D}, m)$. Therefore, carrying out model selection by attempting to minimize these types of generalization bounds is closely related to minimizing $\text{stab}(A, \mathcal{D}, m)$. In general, the generalization bound might converge to 0 as $m \rightarrow \infty$, but this is immaterial for the purpose of model selection. The important factor is the relative values of the measure, over different choices of the algorithm parameters Θ . In other words, the important quantity is the relative convergence rates of this bound for different choices of Θ , governed by $\text{stab}(A, \mathcal{D}, m)$.

This informal discussion only exemplifies the relationship between generalization and stability, since the setting and the definition of $d_{\mathcal{D}}$ here differs from the one we will focus on later in the paper. Although these ideas can be generalized, they go beyond the scope of this paper, and we leave it for future work.

²Where we define $D_{KL}[\mathcal{Q}||\mathcal{P}] = \int_{\mathcal{X}} \mathcal{Q}(X) \ln(\mathcal{Q}(X)/\mathcal{P}(X))$, and $D_{KL}[q||p]$ for $q, p \in [0, 1]$ is defined as the divergence of Bernoulli distributions with parameters q and p .

4 Effective model selection for arbitrarily large sample sizes

From now on, following [2], we will define the clustering distance function $d_{\mathcal{D}}$ of Eq. (1) as:

$$d_{\mathcal{D}}(A(S_1), A(S_2)) = \Pr_{\mathbf{x}_1, \mathbf{x}_2 \sim \mathcal{D}} (A(S_1)(\mathbf{x}_1, \mathbf{x}_2) \neq A(S_2)(\mathbf{x}_1, \mathbf{x}_2)). \quad (2)$$

In other words, the clustering distance is the probability that two independently drawn instances from \mathcal{D} will be in the same cluster under one clustering, and in different clusters under another clustering.

In [2], it is essentially proven that if there exists a unique optimizer to the clustering algorithm's objective function, to which the algorithm converges for asymptotically large samples, then $\text{stab}(A, \mathcal{D}, m)$ converges to 0 as $m \rightarrow \infty$, regardless of the parameters of A . From this, it was concluded that using stability as a tool for cluster validity is problematic, since for large enough samples it would always be approximately zero, for any algorithm parameters chosen.

However, using the intuition gleaned from the results of the previous section, the different *convergence rates* of the stability measure (for different algorithm parameters) should be more important than their absolute values or the sample size. The key technical result needed to substantiate this intuition is the following theorem:

Theorem 2. *Let X, Y be two random variables bounded in $[0, 1]$, and with strictly positive expected values. Assume $\mathbb{E}[X]/\mathbb{E}[Y] \geq 1 + c$ for some positive constant c . Letting X_1, \dots, X_m and Y_1, \dots, Y_m be m identical independent copies of X and Y respectively, define $\hat{X} = \frac{1}{m} \sum_{i=1}^m X_i$ and $\hat{Y} = \frac{1}{m} \sum_{i=1}^m Y_i$. Then it holds that:*

$$\Pr(\hat{X} \leq \hat{Y}) \leq \exp\left(-\frac{1}{8}m\mathbb{E}[X] \left(\frac{c}{1+c}\right)^4\right) + \exp\left(-\frac{1}{4}m\mathbb{E}[X] \left(\frac{c}{1+c}\right)^2\right).$$

The importance of this theorem becomes apparent when \hat{X}, \hat{Y} are taken to be empirical estimators of $\text{stab}(A, \mathcal{D}, m)$ for two different algorithm parameter sets Θ, Θ' . For example, suppose that according to our stability measure (see Eq. (1)), a cluster model with k clusters is more stable than a model with k' clusters, where $k \neq k'$, for sample size m (e.g., $\text{stab}(A_k, \mathcal{D}, m) < \text{stab}(A_{k'}, \mathcal{D}, m)$). These stability measures might be arbitrarily close to zero. Assume that with high probability over the choice of samples S_1 and S_2 of size m , we can show that $d_{\mathcal{D}}(A_k(S_1), A_k(S_2)) \leq 1/\sqrt{m}$, while $d_{\mathcal{D}}(A_{k'}(S_1), A_{k'}(S_2)) \geq 1.01/\sqrt{m}$. We cannot compute these exactly, since the definition of $d_{\mathcal{D}}$ involves an expectation over the unknown distribution \mathcal{D} (see Eq. (2)). However, we can estimate them by drawing another sample S_3 of m instance pairs, and computing a sample mean to estimate Eq. (2). According to Thm. 2, since $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$ and $d_{\mathcal{D}}(A_{k'}(S_1), A_{k'}(S_2))$ have slightly different convergence rates ($c \geq 0.01$), which are slower than $\Theta(1/m)$, then we can discern which number of clusters is more stable, with a high probability which actually *improves* as m increases.

Therefore, we can use Thm. 2 as a guideline for when a stability estimator might be useful for arbitrarily large sample sizes. Namely, we need to show it is an expected value of some random variable, with at least slightly different convergence rates for different model selections, and with at least some of them dominating $\Theta(1/m)$. We would expect these conditions to hold under quite general settings, since most stability measures are based on empirically estimating the mean of some random variable. Moreover, a central-limit theorem argument leads us to expect an asymptotic form of $\Omega(1/\sqrt{m})$, with the exact constants dependent on the model. This convergence rate is slow enough for the theorem to apply. The difficult step, however, is showing that the differing convergence rates can be detected empirically, without knowledge of \mathcal{D} . In the example above, this reduces to showing that with high probability over S_1 and S_2 , $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$ and $d_{\mathcal{D}}(A_{k'}(S_1), A_{k'}(S_2))$ will indeed differ by some constant ratio independent of m .

Proof of Thm. 2. Using a relative entropy variant of Hoeffding's bound [7], we have that for any $1 > b > 0$ and $1/\mathbb{E}[Y] > a > 1$, it holds that:

$$\begin{aligned} \Pr\left(\hat{X} \leq b\mathbb{E}[X]\right) &\leq \exp(-m D_{KL}[b\mathbb{E}[X] \parallel \mathbb{E}[X]]), \\ \Pr\left(\hat{Y} \geq a\mathbb{E}[Y]\right) &\leq \exp(-m D_{KL}[a\mathbb{E}[Y] \parallel \mathbb{E}[Y]]). \end{aligned}$$

By substituting the bound $D_{KL}[p||q] \geq (p - q)^2/2 \max\{p, q\}$ in the two inequalities, we get:

$$\Pr\left(\hat{X} \leq b\mathbb{E}[X]\right) \leq \exp\left(-\frac{1}{2}m\mathbb{E}[X](1-b)^2\right) \quad (3)$$

$$\Pr\left(\hat{Y} \geq a\mathbb{E}[Y]\right) \leq \exp\left(-\frac{1}{2}m\mathbb{E}[Y]\left(a + \frac{1}{a} - 2\right)\right), \quad (4)$$

which hold whenever $1 > b > 0$ and $a > 1$. Let $b = 1 - (1 - \mathbb{E}[Y]/\mathbb{E}[X])^2/2$, and $a = b\mathbb{E}[X]/\mathbb{E}[Y]$. It is easily verified that $b < 1$ and $a > 1$. Substituting these values into the r.h.s of Eq. (3), and to both sides of Eq. (4), and after some algebra, we get:

$$\Pr(\hat{X} \leq b\mathbb{E}[X]) \leq \exp\left(-\frac{1}{8}m\mathbb{E}[X]\left(\frac{c}{1+c}\right)^4\right),$$

$$\Pr(\hat{Y} \geq b\mathbb{E}[X]) \leq \exp\left(-\frac{1}{4}m\mathbb{E}[X]\left(\frac{c}{1+c}\right)^2\right).$$

As a result, by the union bound, we have that $\Pr(\hat{X} \leq \hat{Y})$ is at most the sum of the r.h.s of the last two inequalities, hence proving the theorem. \square

As a proof of concept, we show that for a certain setting, the stability measure used by [2], as defined above, is meaningful for arbitrarily large sample sizes, even when this measure converges to zero for any choice of the required number of clusters. The result is a simple counter-example to the claim that this phenomenon makes cluster stability a problematic tool.

The setting we analyze is a mixture distribution of three well-separated unequal Gaussians in \mathbb{R} , where an empirical estimate of stability, using a centroid-based clustering algorithm, is utilized to discern whether the data contain 2, 3 or 4 clusters. We prove that with high probability, this empirical estimation process will discern $k = 3$ as much more stable than both $k = 2$ and $k = 4$ (by an amount depending on the separation between the Gaussians). The result is robust enough to hold even if in addition one performs normalization procedures to account for the fact that higher number of clusters entail more degrees of freedom for the clustering algorithm (see [9]).

We emphasize that the simplicity of this setting is merely for the sake of analytical convenience. The proof itself relies on a general and intuitive characteristic of what constitutes a 'wrong' model (namely, having cluster

boundaries in areas of high density), rather than any specific feature of this setting. We are currently working on generalizing this result, using a more involved analysis.

In this setting, by the results of [2], $stab(A_k, \mathcal{D}, m)$ will converge to 0 as $m \rightarrow \infty$ for $k = 2, 3, 4$. The next two lemmas, however, show that the stability measure for $k = 3$ (the 'correct' model order) is smaller than the other two, by a substantial ratio independent of m , and that this will be discerned, with high probability, based on the empirical estimates of $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$. The proofs are technical, and appear in the supplementary material to this paper.

Lemma 1. For some $\mu > 0$, let \mathcal{D} be a Gaussian mixture distribution on \mathbb{R} , with density function

$$p(x) = \frac{2}{3\sqrt{2\pi}} \exp\left(-\frac{(x+\mu)^2}{2}\right) + \frac{1}{6\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) + \frac{1}{6\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2}\right).$$

Assume $\mu \gg 1$, so that the Gaussians are well separated. Let A_k be a centroid-based clustering algorithm, which is given a sample and required number of clusters k , and returns a set of k centroids, minimizing the k -means objective function (sum of squared Euclidean distances between each instance and its nearest centroid). Then the following holds, with $o(1)$ signifying factors which converge to 0 as $m \rightarrow \infty$:

$$stab(A_2, \mathcal{D}, m) \geq \frac{1 - o(1)}{7\sqrt{m}} \exp\left(-\frac{\mu^2}{32}\right), \quad stab(A_4, \mathcal{D}, m) \geq \frac{0.4 - o(1)}{\sqrt{m}}$$

$$stab(A_3, \mathcal{D}, m) \leq \frac{1.1 + o(1)}{\sqrt{m}} \exp\left(-\frac{\mu^2}{8}\right).$$

Lemma 2. For the setting described in Lemma 1, it holds that over the draw of independent sample pairs $(S_1, S_2), (S'_1, S'_2), (S''_1, S''_2)$ (each of size m from \mathcal{D}), the ratio between $d_{\mathcal{D}}(A_2(S'_1), A_2(S'_2))$ and $d_{\mathcal{D}}(A_3(S_1), A_3(S_2))$, as well as the ratio between $d_{\mathcal{D}}(A_4(S'_1), A_4(S'_2))$ and $d_{\mathcal{D}}(A_3(S_1), A_3(S_2))$, is larger than 2 with probability of at least:

$$1 - (4 + o(1)) \left(\exp\left(-\frac{\mu^2}{16}\right) + \exp\left(-\frac{\mu^2}{32}\right) \right).$$

It should be noted that the asymptotic notation is merely to get rid of second-order terms, and is not an essential feature. Also, the constants are by no means the tightest possible. With these lemmas, we can prove that a direct estimation of $\text{stab}(A, \mathcal{D}, m)$, based on a random sample, allows us to discern the more stable model with high probability, for arbitrarily large sample sizes.

Theorem 3. For the setting described in Lemma 1, define the following unbiased estimator $\hat{\theta}_{k,4m}$ of $\text{stab}(A_k, \mathcal{D}, m)$: Given a sample of size $4m$, split it randomly into 3 disjoint subsets S_1, S_2, S_3 of size m, m and $2m$ respectively. Estimate $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$ by computing

$$\frac{1}{m} \sum_{x_i, x_{m+i} \in S_3} \mathbf{1}(A_k(S_1)(x_i, x_{m+i}) \neq A_k(S_2)(x_i, x_{m+i})),$$

where (x_1, \dots, x_m) is a random permutation of S_3 , and return this value as an estimate of $\text{stab}(A_k, \mathcal{D}, m)$. If three samples of size $4m$ each are drawn i.i.d from \mathcal{D} , and are used to calculate $\hat{\theta}_{2,4m}, \hat{\theta}_{3,4m}, \hat{\theta}_{4,4m}$, then

$$\Pr\left(\hat{\theta}_{3,4m} \geq \min\{\hat{\theta}_{2,4m}, \hat{\theta}_{4,4m}\}\right) \leq \exp(-\Omega(\mu^2)) + \exp(-\Omega(\sqrt{m})).$$

Proof. Using Lemma 2, we have that:

$$\Pr\left(\frac{\min\{d_{\mathcal{D}}(A_2(S'_1), A_2(S'_2)), d_{\mathcal{D}}(A_4(S''_1), A_4(S''_2))\}}{d_{\mathcal{D}}(A_3(S_1), A_3(S_2))} \leq 2\right) < \exp(-\Omega(\mu^2)). \quad (5)$$

Denoting the event above as B , and assuming it does not occur, we have that the estimators $\hat{\theta}_{2,4m}, \hat{\theta}_{3,4m}, \hat{\theta}_{4,4m}$ are each an empirical average over an additional sample of size m , and the expected value of $\hat{\theta}_{3,4m}$ is at least twice smaller than the expected values of the other two. Moreover, by Lemma 1, the expected value of $d_{\mathcal{D}}(A_3(S_1), A_3(S_2))$ is $\Omega(1/\sqrt{m})$. Invoking Thm. 2, we have that:

$$\Pr\left(\hat{\theta}_{3,4m} \geq \min\{\hat{\theta}_{2,4m}, \hat{\theta}_{4,4m}\} \mid B^c\right) \leq \exp(-\Omega(\sqrt{m})) \quad (6)$$

Combining Eq. (5) and Eq. (6) yield the required result. \square

5 Experiments

In order to further substantiate our analysis above, some experiments were run on synthetic and real world data, with the goal of performing model selection over the number of clusters k . Our first experiment simulated the setting discussed in section 4 (see figure 1). We tested 3 different Gaussian mixture distributions (with $\mu = 5, 7, 8$), and sample sizes m ranging from 2^5 to 2^{22} . For each distribution and sample size, we empirically estimated $\hat{\theta}_2, \hat{\theta}_3$ and $\hat{\theta}_4$ as described in section 4, using the k -means algorithm, and repeated this procedure over 1000 trials. Our results show that although these empirical estimators converge towards zero, their convergence rates differ, with approximately constant ratios between them. Scaling the graphs by \sqrt{m} results in approximately constant and differing stability measures for each μ . Moreover, the error rate does not increase with sample size, and decreases rapidly to negligible size as the Gaussians become more well separated - exactly in line with Thm. 3. Notice that although in the previous section we assumed a large separation between the Gaussians for analytical convenience, good results are obtained even when this separation is quite small.

For the other experiments, we used the stability-based cluster validation algorithm proposed in [9], which was found to compare favorably with similar algorithms, and has the desirable property of

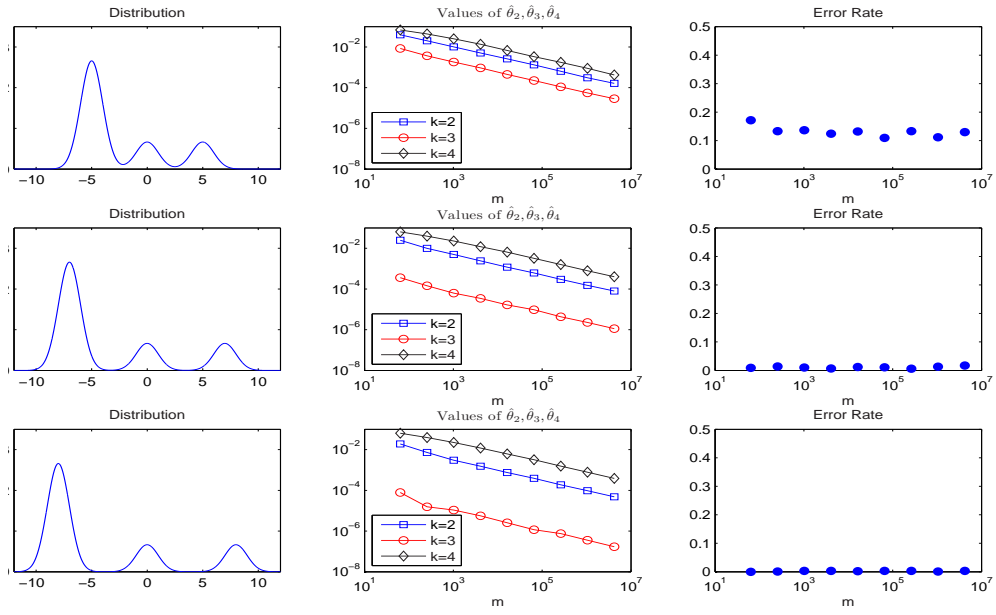


Figure 1: Empirical validation of results in section 4. In each row, the leftmost sub-figure is the actual distribution, the middle sub-figure is a log-log plot of the estimators $\hat{\theta}_2, \hat{\theta}_3, \hat{\theta}_4$ (averaged over 1000 trials), as a function of the sample size, and on the right is the error rate as a function of the sample size (percentage of trials where $\hat{\theta}_3$ was not the smallest of the three).

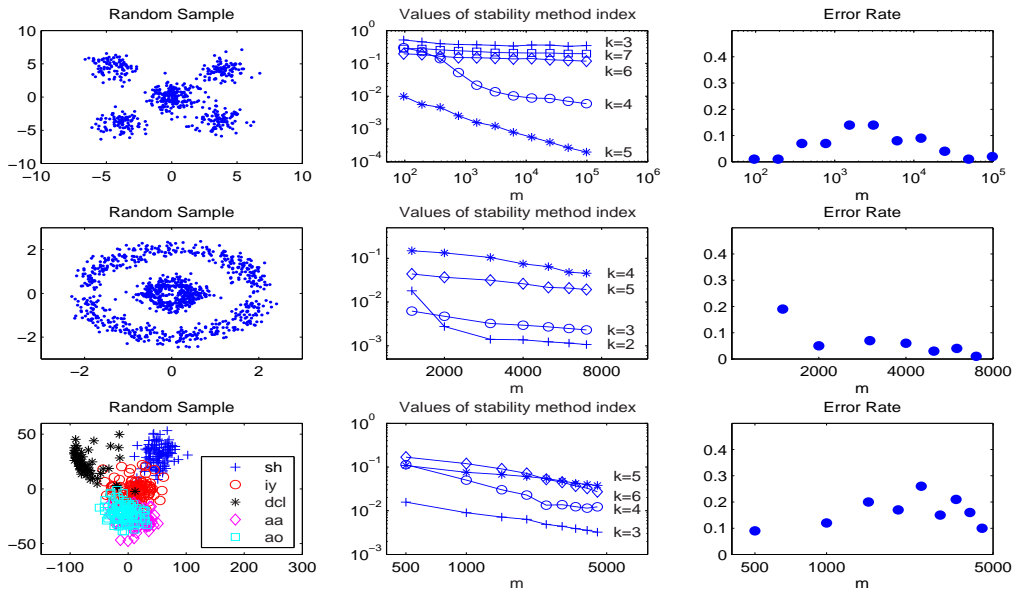


Figure 2: Performance of stability based algorithm in [9] on 3 data sets. For each row, a distribution representation is shown in the left sub-figure; The middle sub-figure is a log-log plot of the computed stability indices (averaged over 100 trials); the right sub-figure is the error rate (percentage of trials where the predicted number of clusters is different from the expected one). In the phoneme data set, the algorithm selects 3 clusters as the most stable models, since the vowels tend to group into a single cluster. The 'errors' are all due to trials when $k = 4$ was deemed more stable.

producing a clear quantitative stability measure, bounded in $[0, 1]$. Lower values match models with higher stability. The synthetic data sets selected (see figure 2) were a mixture of 5 Gaussians, and segmented 2 rings. We also experimented on the Phoneme data set³ [6], which consists of 4,500 log-periodograms of 5 phonemes uttered by English speakers, to which we applied PCA projection on 3 principal components as a pre-processing step. The advantage of this data set is its clear low-dimensional representation relative to its size, allowing us to get nearer to the asymptotic convergence rates of the stability measures. All experiments used the k -means algorithm, except for the ring data set which used the spectral clustering algorithm proposed in [13].

Complementing our theoretical analysis, the experiments clearly demonstrate that regardless of the actual stability measures per fixed sample size, they seem to eventually follow roughly constant and differing convergence rates, with no substantial degradation in performance. In other words, when stability works well for small sample sizes, it should also work at least as well for larger sample sizes. The universal asymptotic convergence to zero does not seem to be a problem in that regard.

6 Conclusions

In this paper, we propose a principled approach for analyzing the utility of stability for cluster validation in large finite samples. This approach stems from viewing stability as a measure of generalization in a statistical setting. It leads us to predict that in contrast to what might be concluded from previous work, cluster stability does not necessarily degrade with increasing sample size. This prediction is substantiated both theoretically and empirically.

The results also provide some guidelines (via Thm. 2) for when a stability measure might be relevant for arbitrarily large sample size, despite asymptotic universal stability. They also suggest that by appropriate scaling, stability measures would become insensitive to the actual sample size used. These guidelines do not presume a specific clustering framework. However, we have proven their fulfillment rigorously only for a certain stability measure and clustering setting. The proof can be generalized in principle, but only at the cost of a more involved analysis. We are currently working on deriving more general theorems on when these guidelines apply.

References

- [1] Shai Ben-David. A framework for statistical clustering with a constant time approximation algorithms for k -median clustering. In *Proceedings of the Seventeenth Annual Conference on Computational Learning Theory*, pages 415–426, 2004.
- [2] Shai Ben-David, Ulrike von Luxburg, and Dávid Pál. A sober look at clustering stability. In *Proceedings of the Nineteenth Annual Conference on Computational Learning Theory*, pages 5–19, 2006.
- [3] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [4] Joachim M. Buhmann and Marcus Held. Model selection in clustering by uniform convergence bounds. In *Advances in Neural Information Processing Systems 12*, pages 216–222, 1999.
- [5] Andrea Caponnetto and Alexander Rakhlin. Stability properties of empirical risk minimization over donsker classes. *Journal of Machine Learning Research*, 6:2565–2583, 2006.
- [6] Trevor Hastie, Robert Tibshirani, Jerome Friedman. *The Elements of Statistical Learning*. Springer, 2001.
- [7] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, March 1963.
- [8] Samuel Kutin and Partha Niyogi. Almost-everywhere algorithmic stability and generalization error. In *Proceeding of the 18th confrence on Uncertainty in Artificial Intelligence (UAI)*, pages 275–282, 2002.
- [9] Tilman Lange, Volker Roth, Mikio L. Braun, and Joachim M. Buhmann. Stability-based validation of clustering solutions. *Neural Computation*, 16(6):1299–1323, June 2004.
- [10] D.A. McAllester. Pac-bayesian stochastic model selection. *Machine Learning Journal*, 51(1):5–21, 2003.
- [11] C. McDiarmid. On the method of bounded differences. In *Surveys in Combinatorics*, volume 141 of *London Mathematical Society Lecture Note Series*, pages 148–188. Cambridge University Press, 1989.
- [12] Alexander Rakhlin and Andrea Caponnetto. Stability of k -means clustering. In *Advances in Neural Information Processing Systems 19*. MIT Press, Cambridge, MA, 2007.
- [13] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [14] Ulrike von Luxburg and Shai Ben-David. Towards a statistical theory of clustering. Technical report, PASCAL workshop on clustering, London, 2005.

³Available at <http://www-stat.stanford.edu/tibs/ElemStatLearn>

Cluster Stability for Finite Samples

Supplementary Material

Ohad Shamir[†] and Naftali Tishby^{†‡}

[†] School of Computer Science and Engineering

[‡] Interdisciplinary Center for Neural Computation

The Hebrew University

Jerusalem 91904, Israel

{ohadsh, tishby}@cs.huji.ac.il

Abstract

This technical report contains the proofs of Lemmas 1 and 2 in the paper 'Cluster Stability for Finite Samples'.

1 Proof of Lemma 1

Proof. The proof idea is essentially identical for all values of k . We have that $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$ is governed by the probability mass of \mathcal{D} which switches between clusters in $A_k(S_1)$ and $A_k(S_2)$, in expectation over S_1 and S_2 . For reasonably large samples, all this probability mass is tightly concentrated in small border regions between the clusters, and is governed by small fluctuations in the border positions. For all k , these fluctuations become smaller as the sample size m increases. The important point is that the location of the border points are different for different choices of k . For the 'right' model, the borders lie in areas of very low probability density, and as a result the probability mass of \mathcal{D} which switches between clusters is relatively small in expectation. In contrast, for the 'wrong' models, some of the border points lie in areas of higher density, so the probability mass of \mathcal{D} which switches between clusters is relatively much higher. From this, we get that $stab(A_k, \mathcal{D}, m)$ is relatively smaller for the 'right' value of k , compared to the other values.

We will consider the case $k = 2$ in some detail, and then go over the other two cases more quickly. To simplify the analysis, the proof involves some approximations, with approximation errors which are asymptotically negligible as $m \rightarrow \infty$, or that are arbitrarily small if μ is large enough. Approximations of the first type form the $o(1)$ term in the lemma, while approximations of the second type can be absorbed into the derived (non-tight) bounds. We will use the formulation $\mathcal{N}(\mu, \sigma^2)$ to denote a normally distributed real random variable, with expectation μ and variance σ^2 . Also, we will make frequent use of the following basic facts: If a_1, a_2 are independent random variables such that $a_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $a_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, then the distribution of $a_1 + a_2$ is $\mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$, and the expected value of $|a_1|$ is $\sqrt{2/\pi}\sigma_1$.

For $k = 2$, let α_1 and α_2 be random variables (over the draw of a sample of size m from \mathcal{D}), representing the centroids in \mathbb{R} returned by the algorithm, such that $\alpha_1 \leq \alpha_2$ (see figure 1). If the Gaussians are well separated, we can assume that they are approximately independent: the value of α_1 is equal to the sample mean derived from the region of the larger Gaussian, while α_2 is equal to the sample mean derived from the mixture of the two smaller Gaussians. The distribution of a sample mean of a unit variance Gaussian is also Gaussian, with variance $1/n$ where n is the sample size on which the mean is estimated. Therefore, we have that the distribution of α_1 is approximately $\mathcal{N}(-\mu, 3/2m)$. Since the two smaller Gaussians are well separated and equal, the distribution of α_2 is approximately the average of the sample means of the Gaussians, namely $\mathcal{N}(\mu/2, 3/m)$.

Let $\beta = (\alpha_1 + \alpha_2)/2$ be a random variable denoting the border point between the two clusters. Since α_1 and α_2 are approximately independent, we have that the distribution of β is approximately

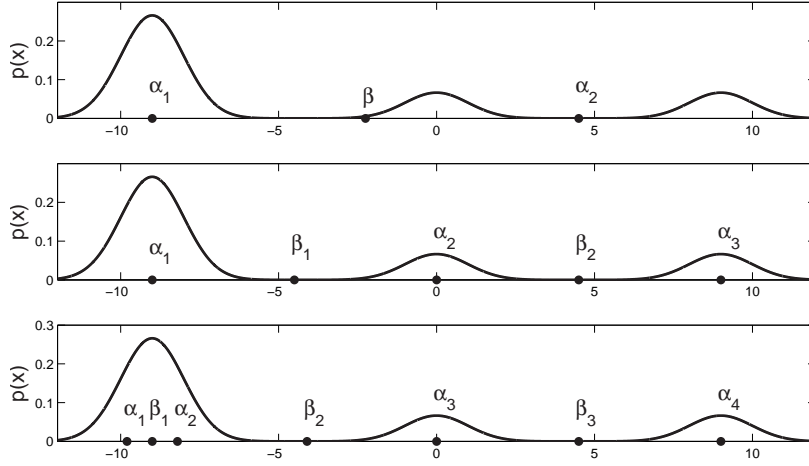


Figure 1: Illustration of centroids and cluster border positions for $k = 2$ (upper sub-figure), $k = 3$ (middle sub-figure), and $k = 4$ (lower sub-figure). The curve represents the density function of \mathcal{D} . For large enough sample sizes, the cluster centroids (denoted by α) and cluster border points (denoted by β) will be tightly concentrated around the positions indicated in the sub-figures.

$\mathcal{N}(-\mu/4, 9/8m)$. As a result, if we let β' and β'' be two independent copies of β , we have that $\beta' - \beta''$ is distributed as $\mathcal{N}(0, 9/4m)$. Finally, since for large values of m we have that β is concentrated around $-\mu/4$, it follows that the probability mass of \mathcal{D} which switches between clusters (over the draw and clustering of two independent samples) is approximately distributed as $|\beta' - \beta''|p(-\mu/4)$, where $p(\cdot)$ is the probability density function of \mathcal{D} . Informally, this is the probability mass which was on 'one side of the border' under the first clustering, and on the 'other side of the border' under the second clustering.

Recall that $d_{\mathcal{D}}(A_k(S_1), A_k(S_2))$ is defined as the probability that two instances sampled from \mathcal{D} will be in the same cluster for clustering $A_k(S_1)$ and in different clusters for clustering $A_k(S_2)$, or vice versa. For $k = 2$ clusters, this reduces to $2t(1 - t)$, where t is a random variable defined over a pair of independent samples S_1 and S_2 , and represents the probability mass of \mathcal{D} which switches clusters between $A_2(S_1)$ and $A_2(S_2)$. By the results of the previous paragraph, t is distributed as $|\beta' - \beta''|p(-\mu/4)$. Therefore, we have that:

$$\begin{aligned}
stab(A_2, \mathcal{D}, m) &= \mathbb{E}[d_{\mathcal{D}}(A_2(S_1), A_2(S_2))] \\
&= \mathbb{E}[2t(1 - t)] \\
&\approx 2\mathbb{E}[p(-\mu/4)|\beta' - \beta''|] - 2\mathbb{E}[(p(-\mu/4))^2(\beta' - \beta'')^2] \\
&\approx \frac{2}{6\sqrt{2\pi}} \exp(-\mu^2/32) \mathbb{E}[|\beta' - \beta''|] - \frac{2}{72\pi} \exp(-\mu^2/16) \text{var}(\beta' - \beta'') \\
&\approx \frac{1}{3\sqrt{2\pi}} \exp(-\mu^2/32) \sqrt{\frac{2}{\pi}} \sqrt{\frac{9}{4m}} - \frac{1}{36\pi} \exp(-\mu^2/16) \frac{9}{4m} \\
&\stackrel{(1)}{\approx} \frac{1}{2\pi\sqrt{m}} \exp(-\mu^2/32) \\
&> \frac{1}{7\sqrt{m}} \exp(-\mu^2/32).
\end{aligned}$$

Step (1) is due to the fact that for large m and/or μ , the second term is negligible compared to the first term.

For $k = 3$ (see figure 1), each centroid is approximately independent and equal to the sample mean of each Gaussian, and therefore the distributions of the two cluster border points β_1 and β_2 are $\mathcal{N}(-\mu/2, 15/8m)$ and $\mathcal{N}(\mu/2, 3/m)$ respectively. Let t_1 denote the probability mass of \mathcal{D} which

switches between the two leftmost clusters (over drawing and clustering two independent samples), and let t_2 denote the probability mass of \mathcal{D} which switches between the two rightmost clusters. Since the two leftmost clusters constitute approximately $5/6$ of the sample, and the two rightmost clusters constitute approximately $1/3$ of the sample, we have that the probability that two instances will be in the same cluster under one clustering, and in different clusters under another clustering, is approximately $2t_1(5/6 - t_1) + 2t_2(1/3 - t_2)$. As before, let β'_1, β''_1 be two identical independent copies of β_1 , and β'_2, β''_2 be two identical independent copies of β_2 . We have that $\beta'_1 - \beta''_1$ is distributed as $\mathcal{N}(0, 15/4m)$ and $\beta'_2 - \beta''_2$ is distributed as $\mathcal{N}(0, 6/m)$. Therefore:

$$\begin{aligned}
stab(A_3, \mathcal{D}, m) &= \mathbb{E}[d_{\mathcal{D}}(A_3(S_1), A_3(S_2))] \\
&\approx \mathbb{E}[2t_1(\frac{5}{6} - t_1)] + \mathbb{E}[2t_2(\frac{1}{3} - t_2)] \\
&\approx \frac{5}{3}\mathbb{E}[t_1] + \frac{2}{3}\mathbb{E}[t_2] \\
&\approx \frac{5}{3}p(-\mu/2)\mathbb{E}[|\beta'_1 - \beta''_1|] + \frac{2}{3}p(\mu/2)\mathbb{E}[|\beta'_2 - \beta''_2|] \\
&\approx \frac{5}{3}\frac{5}{6\sqrt{2\pi}}\exp(-\mu^2/8)\sqrt{\frac{2}{\pi}}\sqrt{\frac{15}{4m}} + \frac{2}{3}\frac{1}{3\sqrt{2\pi}}\exp(-\mu^2/8)\sqrt{\frac{2}{\pi}}\sqrt{\frac{6}{m}} \\
&= \sqrt{\frac{6250}{864\pi^2m}}\exp(-\mu^2/8) + \sqrt{\frac{8}{27\pi^2m}}\exp(-\mu^2/8) \\
&< \frac{1.1}{\sqrt{m}}\exp(-\mu^2/8).
\end{aligned}$$

For $k = 4$ (see figure 1), we have two centroids α_1, α_2 on the larger Gaussian, and two centroids α_3, α_4 on the two smaller Gaussians. In this case, the expected probability mass which switches clusters over different samplings is overwhelmingly in the region between the clusters of α_1 and α_2 , because all other border areas are in low density areas of \mathcal{D} (taking them into account only improves the derived lower bound).

By theorem 2 in [1], the distribution of β_1 has an asymptotically Gaussian distribution, with a variance which for simplicity will be lower bounded by $3/2m^1$.

As a result, if β'_1 and β''_1 are two identical copies of β_1 , we have that $\beta'_1 - \beta''_1$ is approximately distributed as a Gaussian centered on 0 with a variance of at least $3/m$. We can repeat an argument similar to the other cases (and with the same notation) to get that:

$$\begin{aligned}
stab(A_4, \mathcal{D}, m) &= \mathbb{E}[d_{\mathcal{D}}(A_4(S_1), A_4(S_2))] \\
&\geq \mathbb{E}[2t_1(\frac{2}{3} - t_1)] \\
&\approx \frac{4}{3}\mathbb{E}[t_1] \\
&\approx \frac{4}{3}p(-\mu)\mathbb{E}[|\beta'_1 - \beta''_1|] \\
&\geq \frac{8}{3\pi\sqrt{3m}} \\
&> \frac{0.4}{\sqrt{m}}.
\end{aligned}$$

□

¹In fact, this bound on the variance can be derived directly without resorting to the asymptotic assumption. Since β_1 may be viewed as an unbiased estimator of the larger Gaussian's mean, we can get the result by a direct application of the Crámmer-Rao lower bound.

2 Proof of Lemma 2

Proof. $d_{\mathcal{D}}(A_3(S_1), A_3(S_2))$ is a random variable (over the draw of S_1 and S_2). Its expected value is $\text{stab}(A_3, \mathcal{D}, m)$, which by the previous lemma can be upper bounded (up to asymptotically negligible approximation errors) by $1.1 \exp(-\mu^2/8)/\sqrt{m}$. Therefore, by Markov's inequality, we have that

$$\Pr \left(d_{\mathcal{D}}(A_3(S_1), A_3(S_2)) \geq \frac{1}{2\sqrt{m}} \exp(-\mu^2/16) \right) < 2.2 \exp(-\mu^2/16). \quad (1)$$

We now wish to prove a lower bound on $d_{\mathcal{D}}(A_2(S_1), A_2(S_2))$ which would hold with high probability. In the proof of lemma 1, we have shown that the distribution of $d_{\mathcal{D}}(A_2(S_1), A_2(S_2))$ is approximately (up to negligible factors) $2p(-\mu/4)|\beta' - \beta''|$, where $\beta' - \beta''$ has a normal distribution $\mathcal{N}(0, 9/4m)$, and $p(\cdot)$ is the probability density function of \mathcal{D} . Therefore:

$$\begin{aligned} & \Pr \left(d_{\mathcal{D}}(A_2(S_1), A_2(S_2)) < \frac{1}{\sqrt{m}} \exp(-\mu^2/16) \right) \\ & \approx \Pr \left(\frac{1}{3\sqrt{2\pi}} \exp(-\mu^2/32) |\beta' - \beta''| < \frac{1}{\sqrt{m}} \exp(-\mu^2/16) \right) \\ & = \Pr \left(|\beta' - \beta''| < \frac{3\sqrt{2\pi}}{\sqrt{m}} \exp(-\mu^2/32) \right) \\ & \stackrel{(1)}{\approx} 2 \Pr \left(\beta' - \beta'' < \frac{3\sqrt{2\pi}}{\sqrt{m}} \exp(-\mu^2/32) \right) - 1 \\ & \stackrel{(1)}{\approx} \text{erf} \left(2\sqrt{\pi} \exp(-\mu^2/32) \right) \\ & \stackrel{(2)}{\leq} 4 \exp(-\mu^2/32). \end{aligned} \quad (2)$$

Step (1) is by the normal distribution of $\beta' - \beta''$ as specified above, and (2) is due to the bound $\text{erf}(x) \leq 2x/\sqrt{\pi}$ for $x \geq 0$.

In the same way, we can derive a high-probability lower bound on $d_{\mathcal{D}}(A_4(S_1), A_4(S_2))$. In the proof of lemma 1, we have shown that the distribution of $d_{\mathcal{D}}(A_4(S_1), A_4(S_2))$ is approximately (up to negligible factors) $(4/3)p(-\mu)|\beta'_1 - \beta''_1|$, where $\beta'_1 - \beta''_1$ has a normal distribution with variance of at least $3/m$. Repeating the same argument as above, we have that

$$\begin{aligned} & \Pr \left(d_{\mathcal{D}}(A_4(S_1), A_4(S_2)) < \frac{1}{\sqrt{m}} \exp(-\mu^2/16) \right) \\ & \approx \Pr \left(\frac{8}{9\sqrt{2\pi}} |\beta'_1 - \beta''_1| < \frac{1}{\sqrt{m}} \exp(-\mu^2/16) \right) \\ & \approx 2 \Pr \left(\frac{8}{9\sqrt{2\pi}} (\beta'_1 - \beta''_1) < \frac{1}{\sqrt{m}} \exp(-\mu^2/16) \right) - 1 \\ & = 2 \Pr \left(\beta'_1 - \beta''_1 < \frac{9\sqrt{2\pi}}{8\sqrt{m}} \exp(-\mu^2/16) \right) \\ & \leq \text{erf} \left(\frac{3\sqrt{3\pi}}{8} \exp(-\mu^2/16) \right) \\ & \leq \frac{3\sqrt{3}}{4} \exp(-\mu^2/16). \end{aligned} \quad (3)$$

Combining inequalities 1,2,3, using the union bound, and taking into account the approximations along the way, we have that:

$$\begin{aligned} & \Pr \left(\frac{\min \{d_{\mathcal{D}}(A_2(S'_1), A_2(S'_2)), d_{\mathcal{D}}(A_4(S''_1), A_4(S''_2))\}}{d_{\mathcal{D}}(A_3(S_1), A_3(S_2))} \leq 2 \right) \\ & < (4 + o(1)) \left(\exp \left(-\frac{\mu^2}{16} \right) + \exp \left(-\frac{\mu^2}{32} \right) \right) \end{aligned}$$

□

References

- [1] J.A Hartigan. Asymptotic distributions for clustering criteria. *The Annals of Statistics*, 6(1):117–131, 1978.