

Sampling approach to sparse approximation problem: determining degrees of freedom by simulated annealing

Tomoyuki Obuchi and Yoshiyuki Kabashima
 Department of Mathematical and Computing Science,
 Tokyo Institute of Technology, Yokohama 226-8502, Japan

Abstract—The approximation of a high-dimensional vector by a small combination of column vectors selected from a fixed matrix has been actively debated in several different disciplines. In this paper, a sampling approach based on the Monte Carlo method is presented as an efficient solver for such problems. Especially, the use of simulated annealing (SA), a metaheuristic optimization algorithm, for determining degrees of freedom (the number of used columns) by cross validation is focused on and tested. Test on a synthetic model indicates that our SA-based approach can find a nearly optimal solution for the approximation problem and, when combined with the CV framework, it can optimize the generalization ability. Its utility is also confirmed by application to a real-world supernova data set.

I. INTRODUCTION

In a formulation of compressed sensing, a sparse vector $\mathbf{x} = (x_i) \in \mathbb{R}^N$ is recovered from a given measurement vector $\mathbf{y} = (y_\mu) \in \mathbb{R}^M$ ($M < N$) by minimizing the residual sum of squares (RSS) under a sparsity constraint as

$$\hat{\mathbf{x}}(K) = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 \right\} \text{ subj. to } \|\mathbf{x}\|_0 \leq K, \quad (1)$$

where $A = (A_{\mu i}) \in \mathbb{R}^{M \times N}$ and $\|\mathbf{x}\|_0$ denote the measurement matrix and the number of non-zero components in \mathbf{x} (ℓ_0 -norm), respectively [1]. The need to solve similar problems also arises in many contexts of information science such as variable selection in linear regression, data compression, denoising, and machine learning. We hereafter refer to the problem of (1) as the *sparse approximation problem* (SAP).

Despite the simplicity of its expression, SAP is highly nontrivial to solve. Finding the exact solution of (1) has proved to be NP-hard [2], and various approximate methods have been proposed. Orthogonal matching pursuit (OMP) [3], in which the set of used columns is incremented in a greedy manner to minimize RSS, is a representative example of such approximate solvers. Another way of finding an approximate solution of (1) is by converting the problem into a Lagrangian form $\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_p^p$ in conjunction with relaxing $\|\mathbf{x}\|_0$ to $\|\mathbf{x}\|_p^p = \sum_{i=1}^N |x_i|^p$. In particular, setting $p = 1$ makes the converted problem convex and allows us to efficiently find its unique minimum solution. This approach is often termed LASSO [4]. When prior knowledge about the generation of \mathbf{x} and \mathbf{y} is available in the form of probability distributions, one may resort to the Bayesian framework for inferring \mathbf{x} from

\mathbf{y} . This can be efficiently carried out by approximate message passing (AMP) [5].

Solvers of these kinds have their own advantages and disadvantages, and the choice of which to employ depends on the imposed constraints and available resources. This means that developing novel possibilities is important for offering more choices. Supposing a situation where one can make use of relatively high computational resources, we explore the abilities and limitations of another approximate approach, i.e., Monte Carlo (MC)-based sampling.

In an earlier study, we showed that a version of simulated annealing (SA) [6], which is a versatile MC-based metaheuristic for functional optimization, has the ability to efficiently find a nearly optimal solution for SAP in a wide range of system parameters [7]. In this paper, we particularly focus on the problem of determining the *degrees of freedom*, K , by cross validation (CV) utilizing SA. We will show that the necessary computational cost of our algorithm is bounded by $O(M^2 N |S_K| K_{\max})$, where S_K is the set of tested values of K and $|S_K|$ is its cardinality, and K_{\max} is the largest one among tested values of K . Admittedly, this computational cost is not cheap. However, our algorithm is easy to parallelize and we expect that large-scale parallelization will significantly diminish this disadvantage and will make the CV analysis using SA practical.

II. SAMPLING FORMULATION AND SIMULATED ANNEALING

Let us introduce a binary vector $\mathbf{c} = (c_i) \in \{0, 1\}^N$, which indicates the column vectors used to approximate \mathbf{y} : If $c_i = 1$, the i th column of A , \mathbf{a}_i , is used; if $c_i = 0$, it is not used. We call this binary variable *sparse weight*. Given \mathbf{c} , the optimal coefficients, $\mathbf{x}(\mathbf{c})$, are expressed as

$$\mathbf{x}(\mathbf{c}) = \arg \min_{\mathbf{x}} \|\mathbf{y} - A(\mathbf{c} \circ \mathbf{x})\|_2^2, \quad (2)$$

where $(\mathbf{c} \circ \mathbf{x})_i = c_i x_i$ represents the Hadamard product. The components of $\mathbf{x}(\mathbf{c})$ for the zero components of \mathbf{c} are actually indefinite, and we set them to be zero. The corresponding RSS is thus defined by

$$\mathcal{E}(\mathbf{c}|\mathbf{y}, A) = M\epsilon(\mathbf{c}|\mathbf{y}, A) = \frac{1}{2} \|\mathbf{y} - A\mathbf{x}(\mathbf{c})\|_2^2. \quad (3)$$

To perform sampling, we employ the statistical mechanical formulation in Ref. [7] as follows. By regarding \mathcal{E} as an “energy” and introducing an “inverse temperature” β , we can define a Boltzmann distribution as

$$P(\mathbf{c}|\beta; \mathbf{y}, A) = \frac{1}{G} \delta \left(\sum_i c_i - K \right) e^{-\beta \mathcal{E}(\mathbf{c}|\mathbf{y}, A)}, \quad (4)$$

where G is the “partition function”

$$G = G(\beta|\mathbf{y}, A) = \sum_{\mathbf{c}} \delta \left(\sum_i c_i - K \right) e^{-\beta \mathcal{E}(\mathbf{c}|\mathbf{y}, A)}. \quad (5)$$

In the limit of $\beta \rightarrow \infty$, (4) is guaranteed to concentrate on the solution of (1). Therefore, sampling \mathbf{c} at $\beta \gg 1$ offers an approximate solution of (1).

Unfortunately, directly sampling \mathbf{c} from (4) is computationally difficult. To resolve this difficulty, we employ a Markov chain dynamics whose equilibrium distribution accords with (4). Among many choices of such dynamics, we employ the standard Metropolis-Hastings rule [8]. A noteworthy remark is that a trial move of the sparse weights, $\mathbf{c} \rightarrow \mathbf{c}'$, is generated by “pair flipping” two sparse weights, one equal to 0 and the other equal to 1. Namely, choosing an index i of the sparse weight from ONES $\equiv \{k|c_k = 1\}$ and another index j from ZEROS $\equiv \{k|c_k = 0\}$, we set $\mathbf{c}' = \mathbf{c}$, except for the counterpart of $(c_i, c_j) = (1, 0)$, which is given as $(c'_i, c'_j) = (0, 1)$. This flipping can keep the sparsity constant during the update. A pseudo-code of the MC algorithm is given in Algorithm 1 as a reference.

The most time-consuming part in the algorithm is the evaluation of \mathcal{E}' denoted in the sixth line; the naive operation for it requires $O(MK^2 + K^3)$ since matrix inversion of the gram matrix $A(\mathbf{c})^T A(\mathbf{c})$ is involved, where T and $A(\mathbf{c})$ stand for the matrix transpose and the submatrix of A that is composed of column vectors of A whose column indices belong to ONES, respectively. However, since the flip $\mathbf{c} \rightarrow \mathbf{c}'$ changes $A(\mathbf{c})$ only by two columns, one can reduce this computational cost to $O(MK + K^2)$ using a matrix inversion formula [7]. This implies that, when the average number of flips per variable (MC steps) is kept to a constant, the computational cost of the algorithm scales as $O(MNK)$ per MC step in the dominant order since $M > K$.

In general, a longer time is required for equilibrating MC dynamics as β grows larger. Furthermore, the dynamics has the risk of being trapped by trivial local minima of (3) if β is fixed to a very large value from the beginning. A practically useful technique for avoiding these inconveniences is to start with a sufficiently small β and gradually increase it, which is termed simulated annealing (SA) [6]. As $\beta \rightarrow \infty$, \mathbf{c} is no longer updated, and final configuration \mathbf{c}_{fin} is expected to lead to a solution that is very close (or identical) to the optimal solution in (1), i.e., $\hat{\mathbf{x}}(K) \approx \mathbf{x}(\mathbf{c}_{\text{fin}})$.

SA is mathematically guaranteed to find the globally optimal solution of (1) if β is increased to infinity slowly enough in such a way that $\beta(t) < C \log(t+2)$ is satisfied, where t is the counter of the MC dynamics and C is a time-independent

constant [9]. Of course, this schedule is practically meaningless, and a much faster schedule is employed generally. In Ref. [7], we examined the performance of SA with a very rapid annealing schedule for a synthetic model whose properties of fixed β can be analytically evaluated. Comparison between the analytical and the experimental results indicates that the rapid SA performs quite well unless a phase transition of a certain type occurs at relatively low β . Owing to the analysis of the synthetic model, the range of system parameters in which the phase transition occurs is rather limited. We therefore expect that SA serves as a promising approximate solver for (1).

Algorithm 1 MC update with pair flipping

```

1: procedure MCPF( $\mathbf{c}, \beta, \mathbf{y}, A$ ) ▷ MC routine
2:   ONES  $\leftarrow \{k|c_k = 1\}$ , ZEROS  $\leftarrow \{k|c_k = 0\}$ 
3:   randomly choose  $i$  from ONES and  $j$  from ZEROS
4:    $\mathbf{c}' \leftarrow \mathbf{c}$ 
5:    $(c'_i, c'_j) \leftarrow (0, 1)$ 
6:    $(\mathcal{E}, \mathcal{E}') \leftarrow (\mathcal{E}(\mathbf{c}|\mathbf{y}, A), \mathcal{E}(\mathbf{c}'|\mathbf{y}, A))$ 
7:    $p_{\text{accept}} \leftarrow \max(1, e^{-\beta(\mathcal{E}' - \mathcal{E})})$ 
8:   generate a random number  $r \in [0, 1]$ 
9:   if  $r < p_{\text{accept}}$  then
10:      $\mathbf{c} \leftarrow \mathbf{c}'$ 
11:   end if
12:   return  $\mathbf{c}$ 
13: end procedure

```

III. EMPLOYMENT OF SA FOR CROSS VALIDATION

CV is a framework designed to evaluate the generalization ability of statistical models/learning systems based on a given set of data. In particular, we examine the leave-one-out (LOO) CV, but its generalization to the k -fold CV is straightforward.

In accordance with the cost function of (1), we define the generalization error

$$\epsilon_g = \frac{1}{2} \overline{\left(y_{M+1} - \sum_{i=1}^N A_{(M+1),i} x_i \right)^2} \quad (6)$$

as a natural measure for evaluating the generalization ability of \mathbf{x} , where $\overline{\dots}$ denotes the expectation with respect to the “unobserved” $(M+1)$ th data $(\{A_{(M+1),i}\}, y_{M+1})$. LOO CV assesses the LOO CV error (LOOE)

$$\epsilon_{\text{LOO}}(K|\mathbf{y}, A) = \frac{1}{2M} \sum_{\mu=1}^M \left(y_{\mu} - \sum_{i=1}^N A_{\mu i} x_i^{\setminus \mu}(\mathbf{c}^{\setminus \mu}) \right)^2 \quad (7)$$

as an estimator of (6) for the solution of (1), where $\mathbf{x}^{\setminus \mu}(\mathbf{c}^{\setminus \mu}) = (x_i^{\setminus \mu}(\mathbf{c}^{\setminus \mu}))$ is the solution of (1) for the “ μ th LOO system,” which is defined by removing the μ th data $(\{A_{\mu i}\}, y_{\mu})$ from the original system. One can evaluate (7) by independently applying SA to each of the M LOO systems.

The LOOE of (7) depends on K through $\mathbf{c}^{\setminus \mu}$, and hence we can determine its “optimal” value from the minimum of $\epsilon_{\text{LOO}}(K|\mathbf{y}, A)$ by sweeping K . Compared to the case of a single run of SA at a given K , the computational cost for

LOO CV is increased by a certain factor. This factor is roughly evaluated as $O(M \times |S_K| \times K_{\max})$ when varying K in the set of S_K in which the maximum value is K_{\max} . However, this part of the computation can be easily parallelized and, therefore, may not be so problematic when sufficient quantities of CPUs and memories are available. In the next section, the rationality of the proposed approach is examined by application to a synthetic model and a real-world data set.

Before closing this section, we want to remark on two important issues. For this, we assume that \mathbf{y} is generated by a true sparse vector \mathbf{x}_0 as

$$\mathbf{y} = A\mathbf{x}_0 + \boldsymbol{\xi}, \quad (8)$$

where $\boldsymbol{\xi} \in \mathbb{R}^M$ is a noise vector whose entries are uncorrelated with one another.

The first issue is about the accuracy in inferring \mathbf{x}_0 . When A is provided as a column-wisely normalized zero mean random matrix whose entries are uncorrelated with one another, (6) is linearly related to the squared distance between the true and inferred vectors, \mathbf{x}_0 and $\hat{\mathbf{x}}$, as

$$\epsilon_g = \frac{d_1}{N} \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 + d_0 \quad (9)$$

[10]–[12], where d_1 and d_0 are positive constants. Hence, minimizing the estimator of (6), i.e. (7), leads to minimizing the squared distance from the true vector \mathbf{x}_0 . The same conclusion has also been obtained for LASSO in the limit of $M \rightarrow \infty$ while keeping N and K finite [13], where it is not needed to assume absence of correlations in A .

The second issue is about the difficulty in identifying sparse weight vector \mathbf{c}_0 of \mathbf{x}_0 using CV, which is known to fail even when $M/N \rightarrow \infty$ [14]. Actually, we have tried a naive approach to identify \mathbf{c}_0 by directly minimizing a CV error with the use of SA and confirmed that it does not work. A key quantity for this is another type of LOOE:

$$\tilde{\epsilon}_{\text{LOO}}(\mathbf{c}|\mathbf{y}, A) = \frac{1}{2M} \sum_{\mu=1}^M \left(y_{\mu} - \sum_{i=1}^N A_{\mu i} x_i^{\mu}(\mathbf{c}) \right)^2. \quad (10)$$

This looks like (7), but is different in that \mathbf{c} is common among all the terms. It may be natural to expect that the sparse weight minimizing (10), $\tilde{\mathbf{c}} = \text{argmin}_{\mathbf{c}} \{\tilde{\epsilon}_{\text{LOO}}(\mathbf{c}|\mathbf{y}, A)\}$, is the “best” \mathbf{c} that converges to \mathbf{c}_0 in the limit of $M/N \rightarrow \infty$. Unfortunately, this is not true; in fact, $\tilde{\epsilon}_{\text{LOO}}(\mathbf{c}|\mathbf{y}, A)$ in general tends to decrease as K increases, irrespective of the value of $\|\mathbf{x}_0\|_0$ [14]. We have confirmed this by conducting SA, handling $\tilde{\epsilon}_{\text{LOO}}(\mathbf{c}|\mathbf{y}, A)$ as an energy function of \mathbf{c} . Therefore, minimizing (10) can neither identify \mathbf{c}_0 nor offer any clue for determining K .

We emphasize that minimization of (7) can be utilized to determine K to optimize the generalization ability, although it does not have the ability to identify \mathbf{c}_0 either. The difference between these two issues is critical and confusing, as several earlier studies have provided some controversial implications to the usage of CV in sparse inference on linear models [13], [15]–[18].

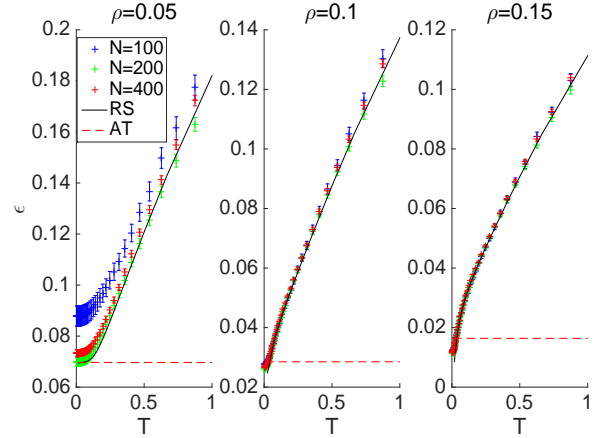


Fig. 1. RSS per component ϵ versus $T = \beta^{-1}$ observed in the annealing process for $N = 100, 200,$ and 400 . Curves represent the RS predictions for (4). The replica symmetry is broken owing to the AT instability below the broken lines.

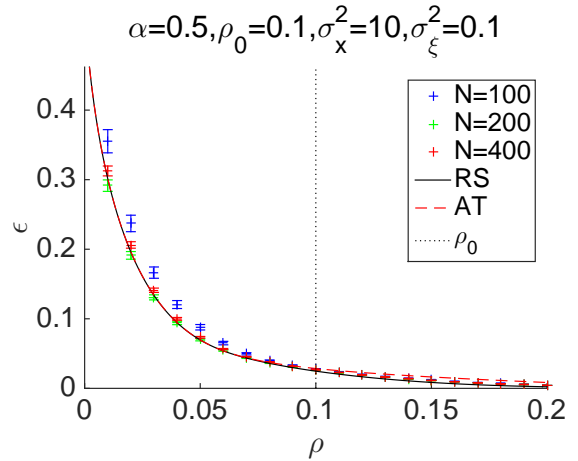


Fig. 2. RSS per component ϵ finally achieved by SA (symbols), its RS assessments for $\beta \rightarrow \infty$ (full curve) and at the onset of the AT instability (red broken curve), plotted against $\rho = K/N$.

IV. RESULT

A. Test on a synthetic model

We first examine the utility of our methodology by applying it to a synthetic model in which vector \mathbf{y} is generated in the manner of (8). For analytical tractability, we assume that A is a simple random matrix whose entries are independently sampled from $\mathcal{N}(0, N^{-1})$ and that each component of \mathbf{x}_0 and $\boldsymbol{\xi}$ is also independently generated from $(1 - \rho_0)\delta(x) + \rho_0\mathcal{N}(0, \sigma_x^2)$ and $\mathcal{N}(0, \sigma_{\xi}^2)$, respectively. Under these assumptions, an analytical technique based on the replica method of statistical mechanics makes it possible to theoretically assess the typical values of various macroscopic quantities when \mathbf{c} is generated from (4) as $N \rightarrow \infty$ keeping $\alpha = M/N$ finite [19]. We performed a theoretical assessment under the so-called *replica symmetric (RS)* assumption.

In the experiment, system parameters were fixed as $\rho_0 =$

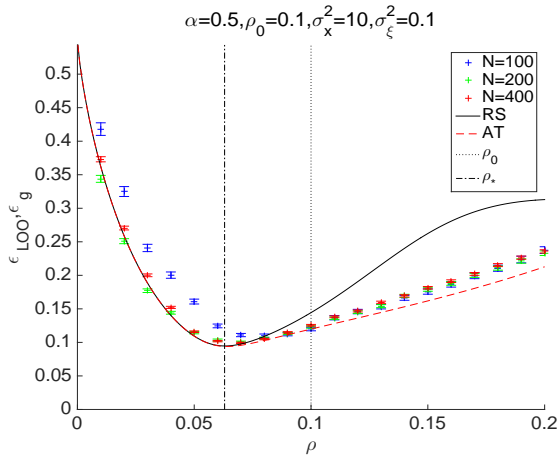


Fig. 3. ϵ_{LOO} evaluated by solutions of SA (symbols) and RS assessments of typical ϵ_g for $\beta \rightarrow \infty$ (full curve) and at the onset of the AT instability (red broken curve), plotted against $\rho = K/N$.

0.1, $\alpha = 0.5$, $\sigma_x^2 = 10$, and $\sigma_\xi^2 = 0.1$. The annealing schedule was set as

$$\beta_a = \beta_0 + r^{a-1} - 1, \quad \tau_a = \tau, \quad (a = 1, \dots, 100), \quad (11)$$

where τ_a denotes the typical number of MC flips per component for a given value of inverse temperature β_a . We set $\tau = 5$, $\beta_0 = 10^{-8}$, and $r = 1.1$ as default parameter values. Thus, the maximum value of β was $\beta_{100} \approx 1.3 \times 10^4$. The examined system sizes were $N = 100, 200$, and 400 . We took the average over $N_{\text{samp}} = 100$ different samples. The error bar is given by the standard deviation among those samples divided by $\sqrt{N_{\text{samp}} - 1}$.

Fig. 1 shows how RSS per component ϵ in (1) depends on $T = \beta^{-1}$ during the annealing process for $K/N \equiv \rho = 0.05, 0.1$, and 0.15 . The data of SA (symbols) for $N = 100, 200$, and 400 totally exhibit a considerably good accordance with the theoretical prediction (curves) for (4) despite the very rapid annealing schedule of (11), in which c_i ($i = 1, 2, \dots, N$) is flipped only five times on average at each value of $\beta = \beta_a$. The replica analysis indicates that the *replica symmetry breaking* (RSB) due to the *de Almeida-Thouless* (AT) instability [20] occurs below the broken lines. However, the RS predictions for $\beta \rightarrow \infty$ still serve as the lower bounds of ϵ_{min} even in such cases [21], [22]. Fig. 2 shows that the values achieved by SA are fairly close to the lower bounds, implying that SA can find nearly optimal solutions of (1).

Let us denote (6) of typical samples generated from (4) at inverse temperature β as $\epsilon_g(\beta)$. The generalization error of the optimal solution of (1) is assessed as $\epsilon_g(\infty)$. Fig. 3 plots the ϵ_{LOO} evaluated by the solutions of SA (symbols) and the RS evaluations of $\epsilon_g(\infty)$ (full curve) and $\epsilon_g(\beta_{\text{AT}})$ (red broken curve), against $\rho = K/N$. Here, β_{AT} is the critical inverse temperature at which RSB occurs owing to the AT instability. The three plots accord fairly well with one another in the left of their minimum point $\rho_* \sim 0.063$, whereas there are considerable discrepancies between $\epsilon_g(\infty)$ and the other two

| K | 1 | 2 | 3 | 4 | 5 |
|-------------------------|--------|--------|--------|--------|--------|
| ϵ_{LOO} | 0.0328 | 0.0239 | 0.0281 | 0.0331 | 0.0334 |

TABLE I
LOO CV ERROR OBTAINED FOR $K = 1-5$ FOR THE TYPE IA SUPERNOVA DATA SET.

| $K = 1$ | | | | | |
|----------------|----|-----|-----|-----|-----|
| variable | 2 | * | * | * | * |
| times selected | 78 | 0 | 0 | 0 | 0 |
| $K = 2$ | | | | | |
| variable | 2 | 275 | * | * | * |
| times selected | 78 | 77 | 1 | 0 | 0 |
| $K = 3$ | | | | | |
| variable | 2 | 1 | 233 | 14 | 69 |
| times selected | 78 | 76 | 69 | 3 | 2 |
| $K = 4$ | | | | | |
| variable | 2 | 1 | 233 | 94 | 225 |
| times selected | 78 | 59 | 56 | 49 | 13 |
| $K = 5$ | | | | | |
| variable | 2 | 36 | 223 | 225 | 6 |
| times selected | 78 | 37 | 33 | 31 | 27 |

TABLE II
THE TOP FIVE VARIABLES SELECTED BY THE $M = 78$ LOO CV FOR $K = 1-5$.

plots for $\rho > \rho_*$.

The discrepancies are considered to be caused by RSB. For $\beta > \beta_{\text{AT}}$, the MC dynamics tends to be trapped by a metastable state. This makes it difficult for SA to find the global minimum of $\epsilon(c)$, which explains why the SA's results are close to $\epsilon_g(\beta_{\text{AT}})$. Fortunately, this trapping works beneficially for the present purpose of raising the generalization ability by lowering ϵ_g , as seen in Fig. 3. As far as we have examined, for fixed ρ , $\epsilon_g(\beta_{\text{AT}})$ never exceeds the RS evaluation of $\epsilon_g(\infty)$ and is always close to ϵ_{LOO} of SA's results. These imply that the generalization ability achieved by tuning $K = N\rho$ using the SA-based CV is no worse than that obtained when CV is performed by exactly solving (1) for LOO systems. This is presumably because, for a large ρ , the optimal solution of (1) overfits the observed (training) data and its generalization ability becomes worse than that of $x(c)$ typically sampled at appropriate values of $\beta (< \infty)$.

B. Application to a real-world data set

We also applied our SA-based analysis to a data set from the SuperNova DataBase provided by the Berkeley Supernova Ia program [23], [24]. Screening based on a certain criteria yields a reduced data set of $M = 78$ and $N = 276$ [25]. The purpose of the data analysis is to select a set of explanatory variables relevant for predicting the absolute magnitude at the maximum of type Ia supernovae by linear regression.

Following a conventional treatment of linear regression, we preprocessed both the absolute magnitude at the maximum (dependent variable) and the 276 candidates of explanatory variables to have zero means. We performed the SA for $M = 78$ LOO systems of the preprocessed data set. The result of

one single experiment with varying K is given in Tables I and II. Table I provides the values of LOOE, which shows that ϵ_{LOO} is minimized at $K = 2$.

Possible statistical correlations between explanatory variables, which were not taken into account in the synthetic model in sec. IV-A, could affect the results of linear regression [26]. The CV analysis also offers a useful clue for checking this risk. Examining the SA results of $M (= 78)$ LOO systems, we could count how many times each explanatory variable was selected, which could be used for evaluating the reliability of the variable [27]. Table II summarizes the results for five variables from the top for $K = 1-5$. This indicates that no variables other than “2,” which stands for *color*, were chosen stably, whereas variable “1,” representing *light curve width*, was selected with high frequencies for $K \leq 3$. Table II shows that the frequency of “1” being selected is significantly reduced for $K \geq 4$. These are presumably due to the strong statistical correlations between “1” and the newly added variables, suggesting the low reliability of the CV results for $K \geq 4$. In addition, for $K \geq 4$, we observed that the results varied depending on samples generated by the MC dynamics in SA, which implies that there exist many local minima in (3) of LOO systems for $K \geq 4$. These observations mean that we could select at most only *color* and *light curve width* as the explanatory variables relevant for the absolute magnitude prediction with a certain confidence. This conclusion is consistent with that of [25], in which the relevant variables were selected by LASSO, combined with hyper parameter determination following the *ad hoc* “one-standard-error rule,” and with the comparison between several resulting models.

V. SUMMARY

We examined the abilities and limitations of simulated annealing (SA) for sparse approximation problem, in particular, when employed for determining degrees of freedom by cross validation (CV). Application to a synthetic model indicates that SA can find nearly optimal solutions for (1), and when combined with the CV framework, it can optimize the generalization ability. Its utility was also tested by application to a real-world supernova data set.

Although we focused on the use of SA, samples at finite temperatures contain useful information for SAP. How to utilize such information is currently under investigation.

ACKNOWLEDGMENTS

This work was supported by JSPS KAKENHI under grant numbers 26870185 (TO) and 25120013 (YK). The UC Berkeley SNDB is acknowledged for permission of using the data set of sec. IV-B. Useful discussions with Makoto Uemura and Shiro Ikeda on the analysis of the data set are also appreciated.

REFERENCES

[1] M. Elad, *Sparse and Redundant Representations*, Springer, 2010.
 [2] B.K. Natarajan, “Sparse approximate solutions to linear systems,” *SIAM Journal on Computing*, vol. 24, pp. 227–234, 1995.

[3] Y.C. Pati, R. Rezaifar and P.S. Krishnaprasad, “Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition,” in *1993 Conference Record of The Twenty-Seventh Asilomar Conference on Signals, Systems and Computers*, 1993, pp. 40–44.
 [4] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, pp. 267–288, 1996.
 [5] F. Krzakala, M. Mézard, F. Sausset, Y. Sun and L. Zdeborová, “Probabilistic reconstruction in compressed sensing: algorithms, phase diagrams, and threshold achieving matrices,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2012, pp. P08009 (1–57), 2012.
 [6] S. Kirkpatrick, C.D. Gelatt Jr and M.P. Vecchi, “Optimization by simulated annealing,” *Science*, vol. 220, pp. 671–680, 1983.
 [7] T. Obuchi and Y. Kabashima, “Sparse approximation problem: how rapid simulated annealing succeeds and fails,” arXiv:1601.01074.
 [8] W.K. Hastings, “Monte Carlo sampling methods using Markov chains and their applications,” *Biometrika*, vol. 57, pp. 97–109, 1970.
 [9] S. Geman and D. Geman, “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 6, pp. 721–741, 1984.
 [10] H.S. Seung, H. Sompolinsky, and N. Tishby, “Statistical mechanics of learning from examples,” *Physical Review A*, vol. 45, pp. 6056–6091, 1992.
 [11] M. Opper and O. Winther, “A mean field algorithm for Bayes learning in large feed-forward neural networks,” in *Advances in Neural Information Processing Systems 9*, 1996, pp. 225–231.
 [12] T. Obuchi and Y. Kabashima, “Cross validation in LASSO and its acceleration,” 2016, arXiv:1601.00881.
 [13] D. Homrighausen and D.J. McDonald, “Leave-one-out cross-validation is risk consistent for lasso,” *Machine Learning*, vol. 97, pp. 65–78, 2014.
 [14] J. Shao, “Linear model selection by cross-validation,” *Journal of the American Statistical Association*, vol. 88, pp. 486–494, 1993.
 [15] O. Bousquet and A. Elisseeff, “Stability and generalization,” *Journal of Machine Learning Research*, vol. 2, pp. 499–526, 2002.
 [16] C. Leng, Y. Lin and G. Wahba, “A note on the lasso and related procedures in model selection,” *Statistica Sinica*, vol. 16, pp. 1273–1284, 2006.
 [17] K. Shalev-Shwartz, O. Shamir, N. Srebro and K. Sridharan, “Learnability and stability in the general learning setting,” in *COLT 2009*, 2009.
 [18] H. Xu and S. Mannor, “Sparse algorithms are not stable: a no-free-lunch theorem,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, pp. 187–193, 2012.
 [19] Y. Nakanishi, T. Obuchi, M. Okada and Y. Kabashima, “Sparse approximation based on a random overcomplete basis,” arXiv:1510.02189.
 [20] J.R.L. de Almeida and D.J. Thouless, “Stability of the Sherrington-Kirkpatrick solution of a spin glass model,” *Journal of Physics A: Mathematical and General*, vol. 11, pp. 983–990, 1978.
 [21] T. Obuchi, K. Takahashi and K. Takeda, “Replica symmetry breaking, complexity, and spin representation in the generalized random energy model,” *Journal of Physics A: Mathematical and Theoretical*, vol. 43, pp. 485004 (1–28), 2010.
 [22] T. Obuchi, “Role of the finite replica analysis in the mean-field theory of spin glasses,” Ph. D thesis submitted to Tokyo Institute of Technology in 2010, arXiv:1510.02189.
 [23] “The UC Berkeley Filippenko Group’s Supernova Database,” <http://heracles.astro.berkeley.edu/sndb/>.
 [24] J.M. Silverman, M. Ganeshalingam, W. Li. and A.V. Filippenko, “Berkeley Supernova Ia Program – III. Spectra near maximum brightness improve the accuracy of derived distances to Type Ia supernovae,” *Mon. Not. R. Astron. Soc.*, vol. 425, pp. 1889–1916, 2012.
 [25] M. Uemura, K.S. Kawabata, S. Ikeda and K. Maeda, “Variable selection for modeling the absolute magnitude at maximum of Type Ia supernovae,” *Publ. Astron. Soc. Japan*, vol. 67, pp. 55 (1–9), 2015.
 [26] D. Belsley, *Conditioning Diagnostics: Collinearity and Weak Data in Regression*, Wiley, 1991.
 [27] N. Meinshausen and P. Bühlmann, “Stability selection,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 72, pp. 417–473, 2010.