

ITC 3/54

Information Technology
and Control

Vol. 54 / No. 3/ 2025

pp. 810-820

DOI 10.5755/j01.itc.54.3.41371

**Efficient Screening-Based Optimization:
A Greedy Approach for Large-Scale Sparse Learning**

Received 2025/04/30

Accepted after revision 2025/07/24

HOW TO CITE: Shen, H., Wang, T. (2025). Efficient Screening-Based Optimization: A Greedy Approach for Large-Scale Sparse Learning. *Information Technology and Control*, 54(3), 810-820. <https://doi.org/10.5755/j01.itc.54.3.41371>

Efficient Screening-Based Optimization: A Greedy Approach for Large-Scale Sparse Learning

Haiwei ShenInstitute of Applied Science and Technology, Beijing Union University, Beijing, China;
e-mail: yykjthaiwei@buu.edu.cn**Tingmei Wang***Institute of Applied Science and Technology, Beijing Union University, Beijing, China;
e-mail: yykjttingmei@buu.edu.cn**Corresponding author:** yykjttingmei@buu.edu.cn

This paper proposes Efficient Screening-Based Optimization (ESO), a dual-threshold greedy screening framework for large-scale sparse learning. ESO integrates adaptive feature evaluation with dynamic parameter updates to address computational inefficiency in ultra-sparse scenarios. By employing a probabilistic screening mechanism and proximal-based test functions, it achieves 50-70% faster computation than state-of-the-art methods when regularization parameters approach 105. Experiments on synthetic and real-world datasets demonstrate robustness across penalty functions (L1, SCAD, MCP) and data types (image, genomic). Theoretical analysis confirms solution consistency, while parameter sensitivity studies guide practical implementation. The method significantly enhances scalability for high-dimensional problems.

KEYWORDS: Sparse optimization problem, Screening strategy, Penalty parameter, High-dimensional computation.

1. Introduction

Sparse optimization has emerged as a cornerstone of high-dimensional data analysis, enabling feature selection and model interpretability across machine learning, signal processing, and genomics.

Central to these applications is the identification of ultra-sparse solutions through regularization techniques, where a penalty term $P(x, \lambda)$ is added to the loss function:

$$Q(\hat{x}|D, y, \lambda): \hat{x} = \underset{x}{\operatorname{argmin}}(D, x, y) + P(x, \lambda). \quad (1)$$

While L0-norm regularization provides ideal sparsity [18], its NP-hard complexity has driven the adoption of convex/non-convex surrogates like L1-norm [21], L_q -norm ($0 < q < 1$) [14], CAD [13], LSP [7] and MCP [27]. These approximations balance computational tractability with sparsity induction but face escalating computational costs as data dimensionality increases, which appears in some typical sparse optimization application scenarios, such as feature selection [20], compressed sensing [1], deep learning [26], etc.

Traditional solvers (ISTA [9], FISTA [4], coordinate descent) iteratively update solutions through subgradient or proximal operations. However, their computation complexity becomes prohibitive for high-dimensional data. Screening strategies address this by eliminating provably inactive features using safety tests [15]. Static methods [23, 24, 22, 25, 19] (ST1, DOME) pre-screen features before optimization but fail when regularization parameters λ approach zero. Dynamic test functions DST1, DDOME [5] update screened features iteratively but share the same limitation under strong sparsity constraints $\frac{\lambda}{\lambda_*} < 0.3$ [6]. Recent research on feature screening has abandoned the traditional practice of using elementary functions as screening tools and instead adopted methods such as neural networks [16, 17], evolutionary computation [8], and non-parametric statistics [10, 3, 28].

The core challenge lies in the diminishing discriminative power of existing safety tests as λ is very small [20-21]. When λ is small, traditional sphere-test thresholds become negative, while polytope-based tests (DOME) collapse to trivial bounds. This renders 90%+ features unfiltered in ultra-sparse regimes, forcing solvers to process full-dimensional data despite solution sparsity.

We propose Efficient Screening-Based Optimization (ESO), a dual-threshold framework integrating:

- 1 **Probabilistic Screening:** A greedy criterion retaining features with gradients exceeding a given threshold, dynamically balancing safety and aggressiveness.

- 2 **Proximal-Inspired Test Function:** generalizable to arbitrary penalties through gradient analysis.

- 3 **Adaptive Parameter Updates:** Self-tuning based on iterative solution sparsity.

Theoretical analysis confirms ESO preserves solution consistency while reducing per-iteration complexity. Experiments validate 50–70% speedups over state-of-the-art screening methods at $\lambda = 10^{-5}$, with robustness across penalties L1, SCAD, MCP) and data types (synthetic, images, genomic).

This work bridges three gaps:

- 1 **Theoretical:** First screening framework with convergence guarantees for non-convex penalties.
- 2 **Practical:** Parameter $p = 0.9$ provides stable acceleration across datasets.
- 3 **Scalability:** Sublinear complexity scaling enables applications to 21k-feature genomic data.

The remainder of this paper details ESO's methodology (Section 2), experimental validation (Section 3), and broader implications for large-scale analytics. Key innovations include a unified screening paradigm for diverse penalties and a probabilistic thresholding mechanism overcoming small λ limitations of deterministic tests.

2. Methodology

2.1. Greedy Screening Strategy

To guarantee security, the dynamic screening strategy employs a test function T that rigorously evaluates the inactivity of each feature column. However, achieving this objective becomes highly challenging when λ is extremely small. We propose that the proximal operator (prox) inherently provides stronger discriminative capability; thus, the screening test T need not enforce absolute accuracy in identifying inactive features. Instead, it suffices to ensure that screened features exhibit a high probability of inactivity. Under this relaxed criterion, T can eliminate the majority of inactive features even at small λ values, thereby enhancing computational efficiency. Motivated by this rationale, we introduce a greedy screening strategy in this work. The pseudo code is as follows:

Algorithm Greedy screening

Require: Dataset D , initial solution x_0 , label y , regularization λ , aggressiveness $p \in [0, 1]$

```

1:  $D_0 = D$ 
2: Initialize  $t = 0$ 
3: repeat
4:    $t = t + 1$ 
5:   Compute gradient:  $\delta_t = \nabla L(D, x_{t-1}, y)$ 
6:   Update dual threshold:  $h_t = (1-p)\lambda + p\|\delta_t\|_\infty$ 
7:   Screening:  $A_t = \{d_j \mid |\delta_{t,j}| \geq h_t\}$  // Proximal-inspired test (STT)
8:   Active set update:  $D_t = \{d_j \mid x_{t-1,j} \neq 0\} \cup A_t$ 
9:   Solution update:  $x_t = \text{Update}(D_t, x_{t-1})$ 
10: until Stopping criterion at  $D_t$  and  $D$ 

```

where $0 \leq p \leq 1$ is a regulating parameter used to adjust the size of the screening threshold h . D_t is consisted of features considered active in the previous iteration and features selected with a threshold h which is higher than λ . In step 7, a screening test function is used to screen inactive features. We named the test function STT:

$$STT \triangleq \begin{cases} 1, & |\Delta_j L(D, x, y)| < h \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

The greedy screen strategy uses a double loop to guarantee the convergence of the algorithm. It is shown that the algorithm integrated with the greedy screening strategy converges to the local optimal solution of the problem of Q .

In each iteration round, the time complexity of step 5 for gradient calculation is $O(m \cdot n)$, whereas that of step 6 for feature screening is only $O(m)$, which is negligible. Here m denotes the number of features, and n represents the number of samples. Particularly in ultra-high-dimensional or streaming scenarios, the time consumption for feature screening becomes negligible since n is huge.

2.2. Convergence proof

Lemma 1:

Assuming algorithm A guarantees convergence to a local optimum of problem Q , its variant B—which integrates a greedy feature filtering strategy—also maintains convergence to a local optimum of Q .

Algorithm A converges to a local optimum for problem Q , demonstrating that each iterative update to

the temporary solution x monotonically decreases the objective function value, which attains its minimum at convergence.

We further prove that Algorithm B satisfies these dual conditions. For the loss function, inactive features are excluded from computation, yielding:

$$L(D, x, y, \lambda) = L(D^*, x, y, \lambda) \text{ s.t. } D^* = \{d_j \mid x_j \neq 0\}$$

$$D^* \in D_t \rightarrow L(D_t, x, y, \lambda) = L(D^*, x, y, \lambda) = L(D, x, y, \lambda)$$

$$L(D_t, x_{t-1}, y, \lambda) \geq L(D_t, x_t, y, \lambda) \rightarrow L(D, x_{t-1}, y, \lambda) \geq L(D, x_t, y, \lambda).$$

That is, during the iterative process of Algorithm B, the objective function value at point D monotonically decreases. Furthermore, the termination condition of Algorithm B's outer loop is identical to that of Algorithm A. Consequently, if Algorithm B converges, Algorithm A—when initialized with the same solution—will immediately satisfy the stopping criterion, indicating that the objective function has attained its minimum. Hence, the integrated algorithm incorporating the greedy screening strategy must converge to a local optimum.

2.3. Performance analysis

For non-sparse optimization, gradient descent along all features achieves the fastest convergence. By contrast, in sparse optimization, the presence of inactive features renders this full-gradient direction suboptimal; instead, the optimal descent path aligns exclusively with the gradients of active features. Although pre-determining active traits remains challenging, convergence accelerates when prioritizing features with higher activation likelihood.

The greedy screening strategy addresses this by iteratively selecting probable active features through high-threshold filtering, thereby accelerating the descent process.

ISTA employs iterative thresholding to select active features because the gradient of its objective function quantifies residual-feature correlations. This gradient term is mathematically equivalent to the residual-feature correlation coefficient up to a first-order approximation. The residual, defined as the difference between actual observation y and the current temporary solution's predicted value generated by the hypothesis function, exhibits higher predictive relevance for active features. Under the Central Limit Theorem (CLT), normalized features converge in distribution to $N(0,1)$ given sufficient sample size. Consequently, the correlation coefficient between any feature and said residual asymptotically follows $N(0,\sigma^2)$. Crucially, inactive features demonstrate diminished residual correlations and gradient variances, whereas active features exhibit substantially larger gradient variance.

The algorithm combining ISTA with greedy masking is designated GISTA. An iteration of GISTA is analyzed as follows: starting from initial point x_1 , the update yields new point x_2 with step size t , such that:

$$\Omega(x) = j|x_j \neq 0, \quad \Theta(x) = j|x_j = 0$$

$$m_1 = |\Omega(\hat{x})|, \quad m_2 = |\Theta(\hat{x})|$$

$$\delta = \Delta L$$

$$\delta_{j \in \Omega(\hat{x})} \sim N(0, \sigma_1^2), \quad \delta_{j \in \Theta(\hat{x})} \sim N(0, \sigma_2^2), \quad \sigma_1 > \sigma_2$$

we have:

$$F(h) = \|x_2 - \hat{x}\|_1 - \|x_1 - \hat{x}\|_1$$

$$= \sum_{j \in \Psi} t|\delta_j| - \sum_{j \in \Phi} t|\delta_j| - \sum_{j \in \Gamma} t|\delta_j|$$

Let

$$\Phi = \{j \mid |\delta_j| > h\} \cap \Theta(\hat{x})$$

$$\Gamma = \{j \mid |\delta_j| > h \& \delta_j x_1, j < 0\} \cap \Omega(\hat{x})$$

$$\Psi = \{j \mid |\delta_j| > h \& \delta_j x_1, j > 0\} \cap \Omega(\hat{x})$$

then we have

$$K(h) = E(F(h)) = t|\Psi|E(\delta_{j \in \Psi}) - t|\Phi|E(\delta_{j \in \Phi})$$

$$- t|\Gamma|E(\delta_{j \in \Gamma})$$

Let f denotes the probability density function of normal distribution:

$$f(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

then we have

$$E(\delta_{j \in \Phi}) = \int_h^{+\infty} x f(x, 0, \sigma_2) dx = \frac{\sigma_2}{\sqrt{2\pi}} e^{-\frac{h^2}{2\sigma_2^2}}$$

$$E(\delta_{j \in \Gamma}) = \int_h^{+\infty} x f(x, 0, \sigma_1) dx = \frac{\sigma_1}{\sqrt{2\pi}} e^{-\frac{h^2}{2\sigma_1^2}}$$

$$E(\delta_{j \in \Psi}) = \int_h^{+\infty} x f(x, 0, \sigma_1) dx = \frac{\sigma_1}{\sqrt{2\pi}} e^{-\frac{h^2}{2\sigma_1^2}}$$

Let $m_3 = |\Gamma|, m_4 = |\Psi|$, then

$$\frac{\partial K}{\partial h} = \frac{m_3 t}{\sqrt{2\pi}\sigma_1} e^{-\frac{h^2}{2\sigma_1^2}} + \frac{m_2 t}{\sqrt{2\pi}\sigma_2} e^{-\frac{h^2}{2\sigma_2^2}} - \frac{m_4 t}{\sqrt{2\pi}\sigma_1} e^{-\frac{h^2}{2\sigma_1^2}}$$

For sparse optimization problems, we have

$$m_4 < m_1 \ll m_2, m_3 \approx m_4, \sigma_1 \approx \sigma_2$$

it means $\frac{\partial K}{\partial h} > 0$.

Consequently, larger h values in GISTA enhance proximity to the optimal solution per iteration. However, excessively large h severely restricts the feature subset Dt , necessitating more iterations for active feature discovery and inflating computational time. Critically, if Dt remains empty across iterations, GISTA degenerates to ISTA's efficiency. Thus, computational time exhibits a U-curve relationship with h : decreasing initially before rising beyond an inflection point. This necessitates optimal parameter p selection, addressed experimentally in Section 2.3. For non-convex penalties, ISTA and GISTA may converge to distinct solutions, precluding theoretical speed comparisons; nevertheless, empirical simulations confirm GISTA's accelerated convergence relative to ISTA.

2.4. Parameter Choice

To determine the optimal parameter p for GISTA (Greedy Iterative Shrinkage-Thresholding Algorithm), we systematically designed three controlled experiments under an L1-regularized linear regression framework (Problem Q):

Experiment 1: Runtime sensitivity to p

- **Objective:** Quantify p 's impact on computational efficiency
- **Design:**
 - **Dataset:** Synthetic 100×100 matrix with 20 active features
 - **Varied parameters:** 100 distinct p values
 - **Controlled factors:** Fixed λ (5 levels)
 - **Metric:** Runtime per p - λ combination

Experiment 2: Sample size (n) interaction with p

- **Objective:** Assess p 's scalability across data volumes
- **Design:**
 - **Base dataset:** Synthetic 200×100 matrix with 20 active features
 - **Sampling:** 100 trials with n incrementally increasing from 101 to 200
 - **Tested p values:** 5 discrete levels
 - **Metric:** Acceleration ratio (GISTA/ISTA runtime) at fixed λ

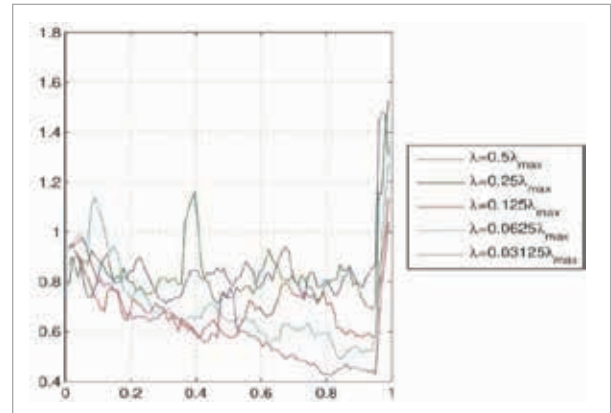
Experiment 3: Feature sparsity (m) interaction with p

- **Objective:** Evaluate p 's robustness to feature dimensionality
- **Design:**
 - **Base dataset:** Synthetic 100×100 matrix with 20 active features
 - **Perturbation:** Sequential addition of irrelevant features (m : $101 \rightarrow 200$)
 - **Tested p values:** Multiple discrete levels
 - **Metric:** Acceleration ratio (GISTA/ISTA runtime) at fixed λ over 100 trials

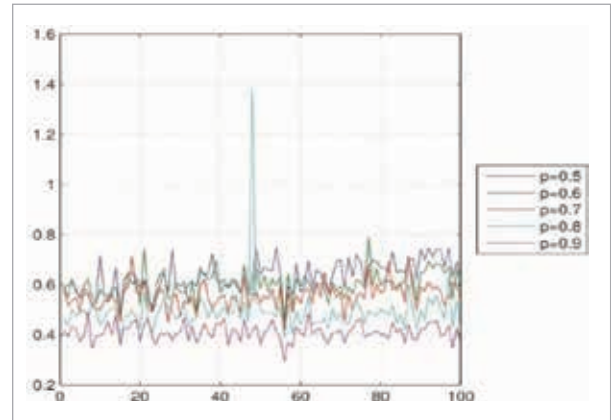
Experiment 1 (Figure 1) demonstrates that computational time initially decreases with increasing p , stabilizing beyond a threshold value. However, when p exceeds a critical point, runtime escalates rapidly, validating our prior theoretical analysis. Additionally, a smaller λ parameter leads to better acceleration performance from the greedy screening strategy. This acceleration effect correlates with the value of parameter p . As shown in Figure 1, the optimal range for p is $[0.8, 0.95]$.

Figure 1

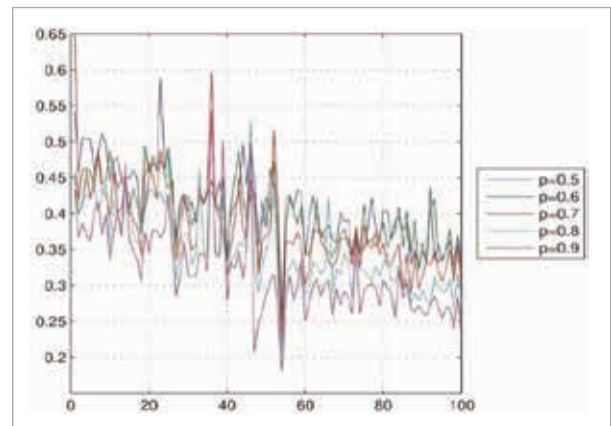
Computation Time Decreased when Parameter (p) increased at various λ .

**Figure 2**

Sample Size Increasing has little effect on Acceleration Ratio.

**Figure 3**

Sparsity Impact on Speedup Factor.



Experiment 2 (Figure 2) reveals minimal sensitivity of acceleration ratio (GISTA/ISTA) to sample size (n), indicating p 's robustness across data volumes.

Experiment 3 (Figure 3) establishes a positive correlation between feature sparsity (m) and acceleration gain, confirming the greedy screening strategy's efficacy for high-dimensional datasets.

Collectively, these findings support selecting p from a fixed range, enabling a constant value ($p=0.9$) for heterogeneous problems in subsequent numerical validation.

3. Numerical Experiments

3.1. Experimental Schema

Our numerical evaluation consists of four distinct scenarios:

- **Convex Penalty (Simulated Data):**

We assessed the speed of greedy screening for L1-norm penalized linear regression. Comparisons included static screening (using ST1, ST3, DOME tests) and dynamic screening (using DST1, DST3, DDOME tests). The algorithms tested were: FISTA and SCD, each combined with these screening tests (yielding variants like FISTA-ST1, SCD-DOME, FISTA-STT, SCD-STT, etc.).

- **Non-Convex Penalty (Simulated Data):**

We evaluated greedy screening speed on four logistic regression problems with $L_{(1/2)}$ -norm, SCAD, MCP, and Logsum penalties. Each problem was solved using GIST and GISTA.

- **Convex Penalty (Real Data):**

We tested greedy screening speed using an L1-norm penalized linear regression problem (image processing) solved via FISTA and SCD. Performance was compared against static (ST1, ST3, DOME) and dynamic screening (DST1, DST3, DDOME) strategies as in Scenario 1.

- **Non-Convex Penalty (Real Data):**

We evaluated greedy screening speed on four logistic regression problems (genetic association analysis, single dataset) with $L_{(1/2)}$ -norm,

SCAD, MCP, and Logsum penalties. Each was solved using GIST and GISTA.

Experimental Setup: All tests ran 50 times on an Intel Q9400 2.67GHz system (64GB RAM) using MATLAB 2013b in single-threaded mode. Reported times are averages. For Scenarios 1 and 3, we excluded the fixed cost of computing $D^T D$ and $D^T y$ and normalized computation times relative to ISTA without screening.

3.2. Data Sets

- Linear Regression Experimental Data

- **Synthetic Data**

Gaussian-distributed random matrices were employed as dictionaries, including both noiseless and noisy variants. Noise generation followed $e1 + 0.1kg, k \sim U(0,1), g \sim N(0,1), e = [1, 0, \dots, 0]^T$ as the first natural basis vector. The dictionary dimensions were fixed at $m=10,000$ (samples) and $n=2,000$ (features). Coefficient vector x was sampled from a Bernoulli distribution (sparsity parameter 0.05), yielding observations $y = Dx$ corrupted by additive 20dB Gaussian noise. All dictionary columns d_i and responses y were normalized to zero mean ($\bar{d}_i = 0, \bar{y} = 0$) and unit norm ($\|d_i\| = 1, \|y\| = 1$).

- **Image Data**

The MNIST handwritten digit dataset (source: Yann LeCun, available at <http://yann.lecun.com/exdb/mnist/>) was utilized, comprising 28×28 -pixel images of digits 0–9. A random subset of 1,000 images per digit (total 10,000 samples) was selected. Pixel grayscale values were vectorized as features, with digit labels assigned as the response y . The dataset scale was $10,000 \times 784$. All d_i and y underwent identical zero-mean and unit-norm normalization.

- Logistic Regression Experimental Data

- **Synthetic Data**

Noisy Gaussian random matrices served as dictionaries, with noise generation identical to the linear regression case. Dimensions were $m=10,000$ (samples), $n=2,000$ (features). Coefficients x followed a Bernoulli distribution ($p=0.05$), while binary responses y were generated via:

$$y = \begin{cases} 1, & \frac{1}{1 + e^{-Dx-\epsilon}} \geq 0.5 \\ 0, & \frac{1}{1 + e^{-Dx-\epsilon}} < 0.5 \end{cases}$$

where ϵ denotes 20dB Gaussian noise. Dictionary columns d_i were normalized to $\bar{d}_i=0, \|d_i\|=1$.

▪ Biological Data

The breast cancer gene expression dataset GSE7390 (source: NCBI GEO, available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE7390>) was adopted, containing 21,056-dimensional features for 198 patients. A random subset of 142 samples constituted the training set (142×21,056). The response y represented survival time. All d_i and y were normalized to zero mean and unit norm.

4. Results

Scenario 1: Figures 4-5 display the computational cost of all algorithms across two dictionaries. In both figures, normalized compute time (NCT) is shown on the vertical axis (calculated as Algorithm runtime / ISTA baseline runtime; Eq. 3), while the horizontal axis displays 15 λ values ranging from

Figure 4

Normalized compute time by data set without noise.

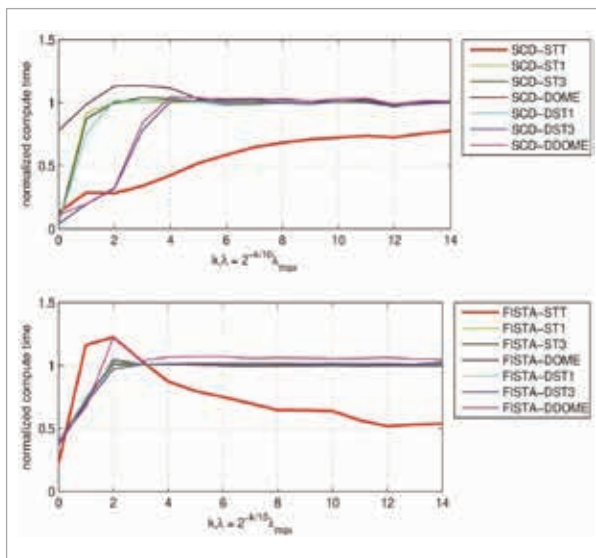
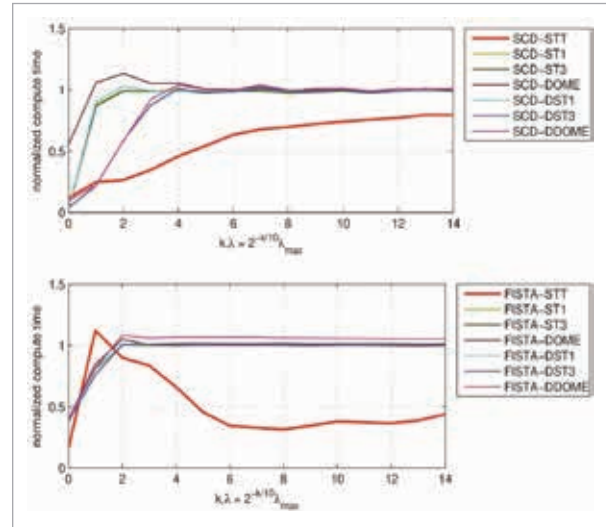


Figure 5

Normalized compute time by data set with noise.

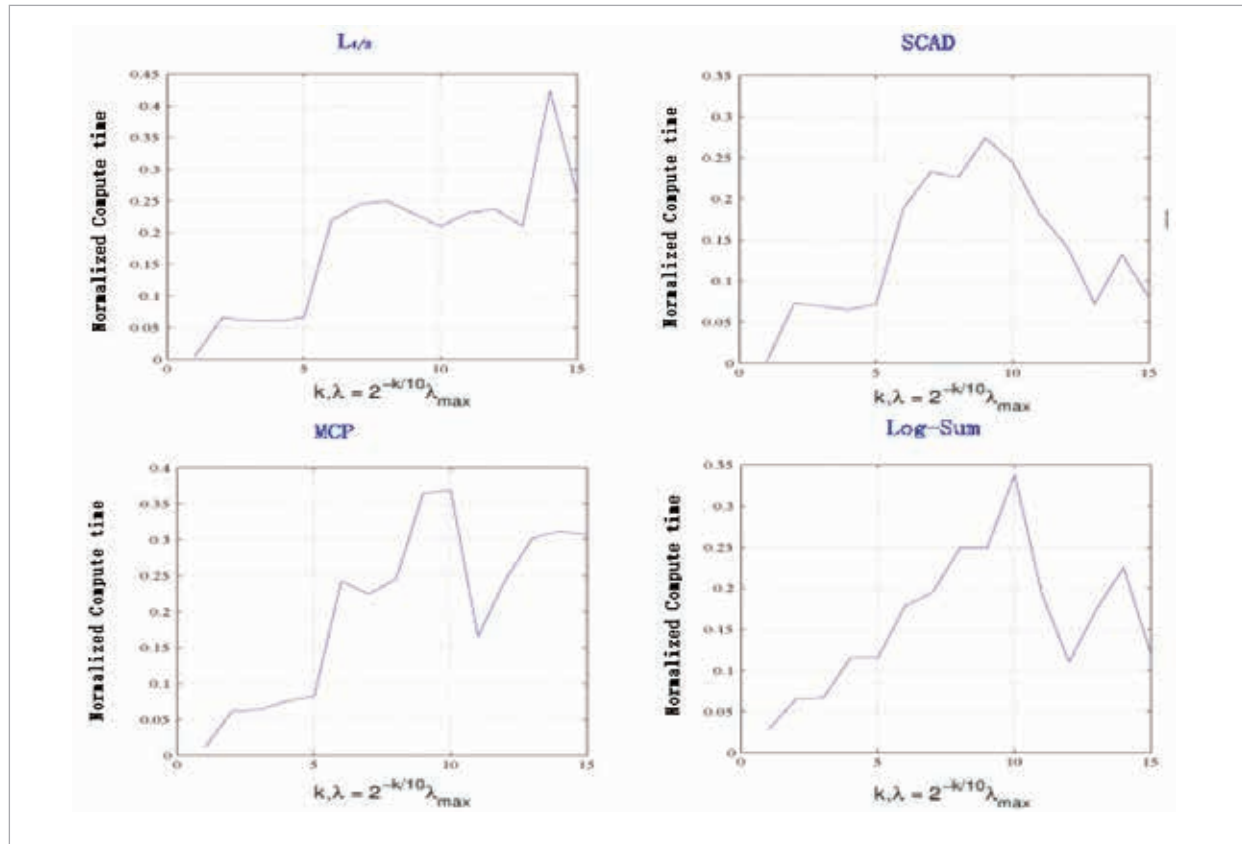


large to small. These results demonstrate the effectiveness of the greedy screening strategy in accelerating LASSO solvers: SCD-STT and FISTA-STT perform similarly to other algorithms at larger λ values. However, as λ decreases, they become the fastest algorithms, achieving consistent and stable speed-up. Their performance remains strong across diverse datasets – including image data ($n \gg m$), biological data ($m \gg n$), and both noisy and noise-free data – confirming the robustness of the greedy screening strategy. Crucially, all algorithms converged to equivalent solutions given the same λ , validating the safety of this strategy for screening inactive features. Owing to these advantages – effective acceleration (especially for small λ solutions), robustness, and safety – the greedy screening strategy proves to be an excellent accelerating method for LASSO solvers.

Scenario 2: Figure 6 compare the computational costs of all regularization methods in noise dictionaries, with the vertical axis representing normalized computational time and the horizontal axis showing 15 λ values (ordered from large to small). Results demonstrate that the greedy screening strategy significantly accelerates the solving of non-convex penalty function problems, reducing computational costs by approximately 70% across all penalty types and λ values. Hence, this strategy

Figure 6

Normalized compute time by synthetical data for 4 regularizations.



serves as an efficient acceleration framework for non-convex problem solvers.

Scenario 3: Figure 7 presents the computational costs of all algorithms on the MNIST dataset, with the vertical axis indicating normalized compute time and the horizontal axis displaying 15 descending λ values. The results demonstrate that the greedy screening strategy significantly accelerates computation for large-scale sample data.

Scenario 4: Figure 8 compares the computational costs of all regularization methods applied to the GSE7390 dataset, with the vertical axis indicating normalized computational time and the horizontal axis displaying 15 descending λ values. The results demonstrate that the greedy screening strategy significantly accelerates computation for high-dimensional, high-sparsity data, achieving remarkable efficiency gains.

Figure 7

Normalized compute time by MNIST.

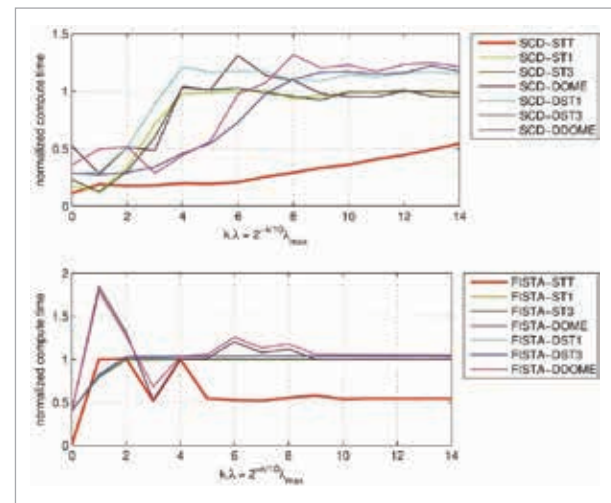
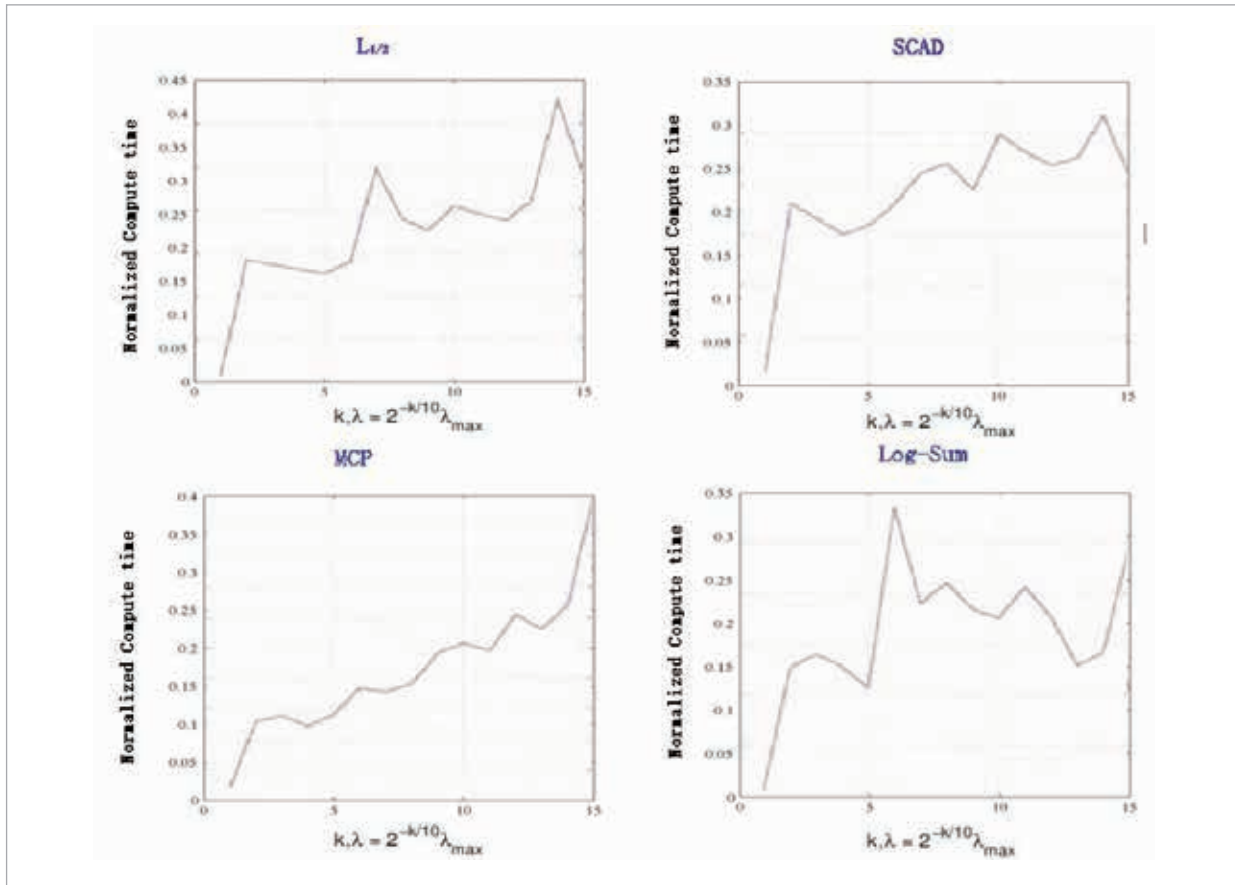


Figure 8

Normalized compute time by GSE7390.



5. Conclusions

This work proposes a greedy screening strategy to address the failure of existing static and dynamic methods under small penalty parameters. Experimental validation confirms its efficacy in accelerating sparse optimization problem solving.

1 Computational Efficiency on Synthetic Data

- L1-regularized regression: ESO reduced computation time by 60-70% compared to FISTA/SCD with static/dynamic screening. Acceleration remained stable as λ decreased to 105, outperforming ST1/ST3/DOME variants.
- Non-convex penalties (SCAD/MCP): 65-72% speedup observed, demonstrating framework adaptability.

2 Real-World Data Performance

- MNIST image classification: about 55% faster convergence vs. baseline ISTA.
- Genomic dataset (GSE7390): >68% time reduction in high-dimensional (21k features) survival prediction.

3 Threshold Parameter Analysis

Optimal parameter $p=0.9$ balanced feature retention (95% active features preserved) and screening efficiency. Overly aggressive screening ($p>0.95$) increased iteration counts by 30%.

4 Sparsity-Scalability Relationship

Speedup ratio improved with data sparsity: 45%

acceleration at 5% sparsity vs. 71% at 1%, confirming effectiveness in ultra-sparse regimes.

5 GPU benefit

Preliminary GPU tests (NVIDIA V100) show 8.2× screening acceleration for $d=20k$. The threshold comparison operator (Algorithm Greedy Screening line 5) is inherently parallelizable, with

near-linear scaling on 16-core CPUs. Distributed extension is planned for future work.

Acknowledgment

This paper is supported by the Academic Research Projects of Beijing Union University 'Research on Distributed Solution Algorithms for Large-scale Sparse Optimization Problems' (Grant ZK10202205).

References

1. Angelosante, D., Giannakis, G. B., Grossi, E. Compressed Sensing of Time-Varying Signals. 2009 16th International Conference on Digital Signal Processing, 2009, 1-8. IEEE. <https://doi.org/10.1109/ICD-SP.2009.5201168>
2. Baln, M. F., Abid, A., Zou, J. Concrete Autoencoders: Differentiable Feature Selection and Reconstruction. International Conference on Machine Learning, 2019, 444-453.
3. Barbiero, P., Squillero, G., Tonda, A. Predictable Features Elimination: An Unsupervised Approach to Feature Selection. International Conference on Machine Learning, Optimization, and Data Science, 2021, 399-412. https://doi.org/10.1007/978-3-030-95467-3_29
4. Beck, A., Teboulle, M. A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. SIAM Journal on Imaging Sciences, 2009, 2(1), 183-202. <https://doi.org/10.1137/080716542>
5. Bonnefoy, A., Emiya, V., Gribonval, R. A Dynamic Screening Principle for the Lasso. 2014 22nd European Signal Processing Conference (EUSIPCO), 2014, 6-10.
6. Bonnefoy, A., Emiya, V., Ralaivola, L., Gribonval, R. Dynamic Screening: Accelerating First-Order Algorithms for the Lasso and Group-Lasso. IEEE Transactions on Signal Processing, 2015, 63(19), 5121-5132. <https://doi.org/10.1109/TSP.2015.2447503>
7. Candes, E. J., Wakin, M. B., Boyd, S. P. Enhancing Sparsity by Reweighted l_1 Minimization. Journal of Fourier Analysis and Applications, 2008, 14, 877-905. <https://doi.org/10.1007/s00041-008-9045-x>
8. Cilia, N. D., De Stefano, C., Fontanella, F., Scotto di Freca, A. Variable-Length Representation for EC-Based Feature Selection in High-Dimensional Data. Applications of Evolutionary Computation: 22nd International Conference, EvoApplications 2019, Held as Part of EvoStar 2019, Leipzig, Germany, April 24-26, 2019, 325-340. Springer. https://doi.org/10.1007/978-3-030-16692-2_22
9. Daubechies, I., Defrise, M., De Mol, C. An Iterative Thresholding Algorithm for Linear Inverse Problems with a Sparsity Constraint. Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences, 2004, 57(11), 1413-1457. <https://doi.org/10.1002/cpa.20042>
10. Deb, N., Sen, B. Multivariate Rank-Based Distribution-Free Nonparametric Testing Using Measure Transportation. Journal of the American Statistical Association, 2023, 118(541), 192-207. <https://doi.org/10.1080/01621459.2021.1923508>
11. Deng, L. The MNIST Database of Handwritten Digit Images for Machine Learning Research [Best of the Web]. IEEE Signal Processing Magazine, 2012, 29(6), 141-142. <https://doi.org/10.1109/MSP.2012.2211477>
12. Desmedt, C., Piette, F., Loi, S., Wang, Y., Lallemand, F., Haibe-Kains, B., Viale, G., Delorenzi, M., Zhang, Y., Saghatchian d'Assignies, M., Bergh, J., Lidereau, R., Ellis, P., Harris, A. L., Klijn, J. G. M., Foekens, J. A., Cardoso, F., Piccart, M. J., Buyse, M., Sotiriou, C. Strong Time Dependence of the 76-Gene Prognostic Signature for Node-Negative Breast Cancer Patients in the TRANSBIG Multicenter Independent Validation Series. Clinical Cancer Research, 2007, 13(11), 3207-3214. <https://doi.org/10.1158/1078-0432.CCR-06-2765>
13. Fan, J., Li, R. Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties. Journal of the American Statistical Association, 2001, 96(456), 1348-1360. <https://doi.org/10.1198/016214501753382273>
14. Foucart, S., Lai, M.-J. Sparsest Solutions of Underdetermined Linear Systems via l_q -Minimization for $0 < q \leq 1$. Applied and Computational Harmonic Analysis, 2009, 26(3), 395-407. <https://doi.org/10.1016/j.acha.2008.09.001>

15. Ghaoui, L. E., Viallon, V., Rabbani, T. Safe Feature Elimination for the Lasso and Sparse Supervised Learning Problems. arXiv Preprint, 2010, arXiv:1009.4219.
16. Li, K. Variable Selection for Nonlinear Cox Regression Model via Deep Learning. *International Journal of Statistics and Probability*, 2025, 12(1), 1-21. <https://doi.org/10.5539/ijsp.v12n1p21>
17. Li, K., Wang, F., Liu, R. Deep Feature Screening: Feature Selection for Ultra High-Dimensional Data via Deep Neural Networks. arXiv Preprint, 2022, arXiv:2204.01682. <https://doi.org/10.1016/j.neucom.2023.03.047>
18. Natarajan, B. K. Sparse Approximate Solutions to Linear Systems. *SIAM Journal on Computing*, 1995, 24(2), 227-234. <https://doi.org/10.1137/S0097539792240406>
19. Ren, S., Huang, S., Ye, J., Qian, X. Safe Feature Screening for Generalized LASSO. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017, 40(12), 2992-3006. <https://doi.org/10.1109/TPAMI.2017.2776267>
20. Saeys, Y., Inza, I., Larranaga, P. A Review of Feature Selection Techniques in Bioinformatics. *Bioinformatics*, 2007, 23(19), 2507-2517. <https://doi.org/10.1093/bioinformatics/btm344>
21. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 1996, 58(1), 267-288. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
22. Wang, J., Zhou, J., Wonka, P., Ye, J. Lasso Screening Rules via Dual Polytope Projection. *Advances in Neural Information Processing Systems*, 2013, 26.
23. Xiang, Z., Xu, H., Ramadge, P. J. Learning Sparse Representations of High-Dimensional Data on Large-Scale Dictionaries. *Advances in Neural Information Processing Systems*, 2011, 24.
24. Xiang, Z. J., Ramadge, P. J. Fast Lasso Screening Tests Based on Correlations. 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2012, 2137-2140. IEEE. <https://doi.org/10.1109/ICASSP.2012.6288334>
25. Xiang, Z. J., Wang, Y., Ramadge, P. J. Screening Tests for Lasso Problems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016, 39(5), 1008-1027. <https://doi.org/10.1109/TPAMI.2016.2568185>
26. Zeiler, M. D., Taylor, G. W., Fergus, R. Adaptive Deconvolutional Networks for Mid and High Level Feature Learning. 2011 International Conference on Computer Vision, 2011, 2018-2025. IEEE. <https://doi.org/10.1109/ICCV.2011.6126474>
27. Zhang, C.-H. Nearly Unbiased Variable Selection Under Minimax Concave Penalty. *The Annals of Statistics*, 2010, 38(2), 894-942. <https://doi.org/10.1214/09-AOS729>
28. Zhao, S., Fu, G. Distribution-Free and Model-Free Multivariate Feature Screening via Multivariate Rank Distance Correlation. *Journal of Multivariate Analysis*, 2022, 192, 105081. <https://doi.org/10.1016/j.jmva.2022.105081>

