

Solve Lasso

Machine Learning Course

Spring 2015

Tsinghua University

Goal

- Compare several methods for solving Lasso.

Lasso

- Problem : Lasso with no structure among covariates

$$\min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1, \quad \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{X} \in \mathbb{R}^{n \times p}$$

- Each row of \mathbf{X} is a data vector.
- Each column of \mathbf{X} represents a feature.
- \mathbf{y} represents the output, which is a sparse construction from \mathbf{X} . \mathbf{w} is the construction weight.

Let $f(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ be the smooth part, $g(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$

be the non-smooth part, then the former problem can be written as

$$\min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w}) + g(\mathbf{w})$$

Lasso

- Implement the first order methods and compare them
- Solvers
 - SubGradient Descent (SGD), **Proximal Methods (ISTA, FISTA)**, **Coordinate Descent (CD)**, Quadratic Programming (QP), Cone Programming (CP), Reweighted L2 (Re-L2), Least Angle Regression (LARS), etc.

Lasso – ISTA solver

For iteration $t=1$ to \dots

-Step1: Calculate the gradient of the smooth part

$$\nabla_{\omega^t} f(\omega^t) = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\omega^t)$$

-Step2: Approximate f around the current point

$$f(\omega) \approx f(\omega^t) + \nabla_{\omega^t} f(\omega^t)^\top (\omega - \omega^t) + \frac{L}{2} \|\omega - \omega^t\|_2^2$$

-Step 2.1 : Solve the problem with fixed L

$$\omega_L^* = \text{Prox}_{\frac{\lambda}{L}}(\omega^t - \frac{1}{L} \nabla_{\omega^t} f(\omega^t)), \text{ where } \text{Prox}_k(x) = \text{sgn}(x) \max(|x|, k)$$

-Step 2.2 : If L satisfies

$$f(\omega_L^*) \leq f(\omega^t) + \nabla_{\omega^t}^\top (\omega_L^* - \omega^t) + \frac{L}{2} \|\omega_L^* - \omega^t\|_2^2, \text{ set } \omega^{t+1} = \omega_L^*$$

and go to next iteration; else increase L by a factor and go to Step 2.1

Lasso – CD solver

Regular version

For iteration $t=1$ to \dots

For each covariate(feature) j , do

$$\omega_j^{t+1} = \text{Prox}_{\frac{\lambda}{2\mathbf{X}_{\cdot j}^\top \mathbf{X}_{\cdot j}}} \left(\omega_j^t - \frac{\mathbf{X}_{\cdot j}^\top \mathbf{X} \boldsymbol{\omega} - \mathbf{X}_{\cdot j}^\top \mathbf{y}}{\mathbf{X}_{\cdot j}^\top \mathbf{X}_{\cdot j}} \right)$$

Stochastic version

For iteration $t=1$ to \dots

Randomly pick a covariate j , do

$$\omega_j^{t+1} = \text{Prox}_{\frac{\lambda}{2\mathbf{X}_{\cdot j}^\top \mathbf{X}_{\cdot j}}} \left(\omega_j^t - \frac{\mathbf{X}_{\cdot j}^\top \mathbf{X} \boldsymbol{\omega} - \mathbf{X}_{\cdot j}^\top \mathbf{y}}{\mathbf{X}_{\cdot j}^\top \mathbf{X}_{\cdot j}} \right)$$

Extension: Lasso – SGD solver

For iteration $t=1$ to \dots

-Step1: Calculate the (sub)gradient

$$\nabla_{\omega^t}(f(\omega^t) + g(\omega^t)) = -2\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\omega^t) + \lambda \sum_i \text{sgn}(\omega_i^t)$$

-Step2 : Decide step size α

$$\alpha^t = a/(t + b), \quad a, b \text{ are pre-defined constant.}$$

-Step3 : Do gradient descent

$$\omega^{t+1} = \omega^t - \alpha^t \nabla_{\omega^t}(f(\omega^t) + g(\omega^t))$$

Extension: Group Lasso

- Problem : Lasso with group structure among covariates (just solve the l1/l2 case)

$$\min_{\omega \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\omega\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\omega_g\|_2, \quad \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{X} \in \mathbb{R}^{n \times p}$$

Also for simplicity, let $f(\omega) = \|\mathbf{y} - \mathbf{X}\omega\|_2^2$, $h(\omega) = \lambda \sum_{g \in \mathcal{G}} \|\omega_g\|_2$.

\mathcal{G} is a partition over the whole index set.

e.g., when $p = 10$, \mathcal{G} can be $\{\{1, 2, 3\}, \{4, 5\}, \{6, 7, 8, 9, 10\}\}$

- For simplicity, we only consider balanced partitions (i.e, the size of each element in g is equal).
- Compare over subgradient method (SGD), proximal method (ISTA), and block coordinate descent method (BCD)

Group Lasso – SGD solver

- Calculate the subgradient and iteratively do gradient descent as the Lasso case.

Group Lasso – ISTA solver

We Just need to change Step 2.1, comparing to Lasso.

-Step 2.1 : Solve subproblems for each group g separately as:

For each group $g \in \mathcal{G}$, do

$$\omega_{g,L}^* = \text{Prox}_{\frac{\lambda}{L\|\omega\|_2}} \left(\omega_g^t - \frac{1}{L} \nabla_{\omega_g^t} f(\omega_g^t) \right),$$

Group Lasso – BCD solver

Regular version

For iteration $t=1$ to \dots

For each group g , iteratively solve the subproblems using SGD
each decent step i :

Norm of this vector

$$\omega_g^{(t+1)i} = \text{Prox}_{\frac{\lambda \|\cdot\|_2}{L}} \left(\omega_g^{(t+1)i-1} - \frac{1}{L} \nabla_{\omega_g^{(t+1)i-1}} f(\omega) \right)$$

Stochastic version

For iteration $t=1$ to \dots

Randomly pick a group g , solve the subproblems using SGD
each decent step i :

$$\omega_g^{(t+1)i} = \text{Prox}_{\frac{\lambda \|\cdot\|_2}{L}} \left(\omega_g^{(t+1)i-1} - \frac{1}{L} \nabla_{\omega_g^{(t+1)i-1}} f(\omega) \right)$$

Task 1

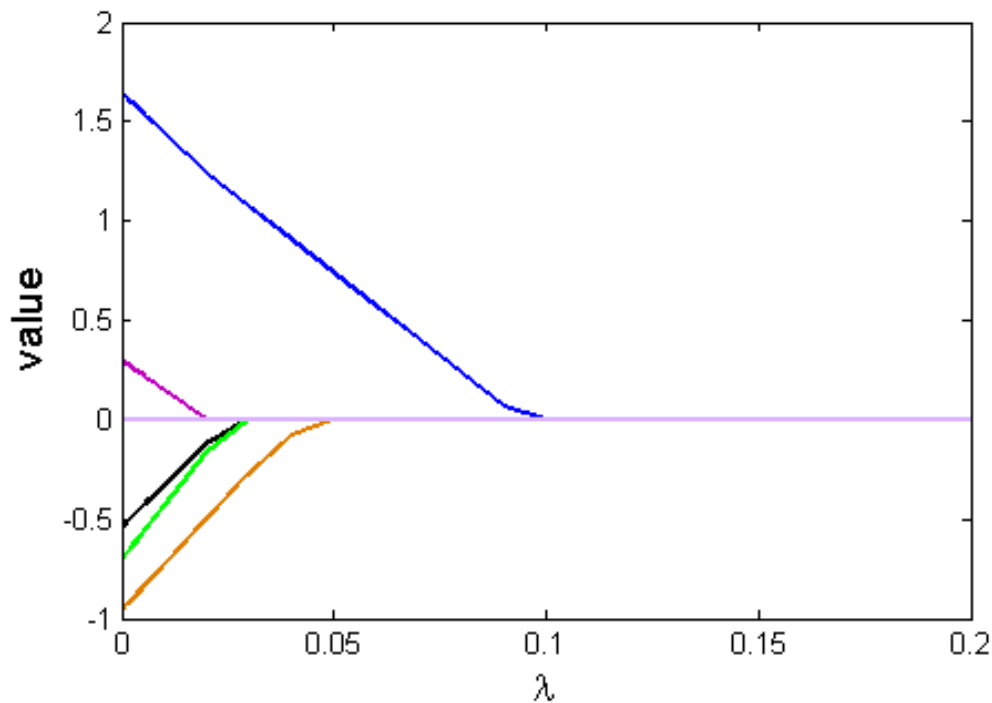
Regular task1 : Draw regularization paths.

E.g., let $n=50$, $p=10$.

1. Draw $\mathbf{X} \in \mathbb{R}^{n \times p}$, where each element was drawn from $\mathcal{N}(0, \frac{1}{n})$
2. Draw $\boldsymbol{\omega}$ whose sparsity level is $s = 0.5$ where each non-zero element was drawn from $\mathcal{N}(0, 1)$
3. Draw noise vector \mathbf{m} from i.i.d. Gaussian $\mathcal{N}(0, 0.01 \|\mathbf{X}\boldsymbol{\omega}\|_2^2/n)$
4. Calculate $\mathbf{y} = \mathbf{X}\boldsymbol{\omega} + \mathbf{m}$
5. Solve Lasso on page 3 (any solver is okay).
6. Draw regularization path (variation of each dimension in $\boldsymbol{\omega}$ when tuning λ).

Task1 : Example

The main point here is the piecewise linearity of paths.
Noise is low in this case so you may get perfect result.



Task 2.1

Regular task 2.1 : Solve Lasso with no correlation among features

Let sparsity level $s = \frac{\# \text{ of zeros in } \boldsymbol{\omega}}{\text{dimension of } \boldsymbol{\omega}}$

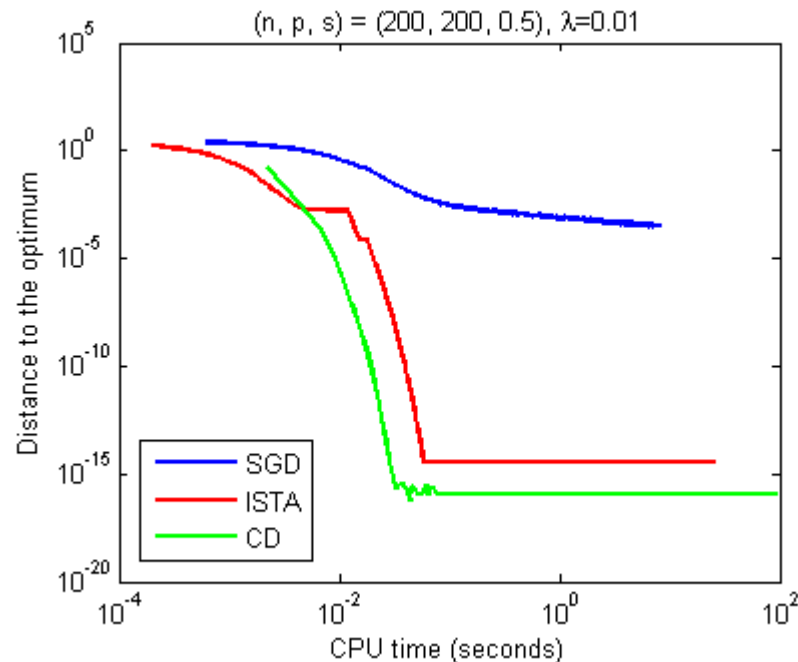
1. Draw $\mathbf{X} \in \mathbb{R}^{n \times p}$, where each element was drawn from $\mathcal{N}(0, \frac{1}{n})$
2. Draw $\boldsymbol{\omega}$ whose sparsity level is s where each non-zero element was drawn from $\mathcal{N}(0, 1)$
3. Draw noise vector \mathbf{m} from i.i.d. Gaussian $\mathcal{N}(0, 0.01 \|\mathbf{X}\boldsymbol{\omega}\|_2^2 / n)$
4. Calculate $\mathbf{y} = \mathbf{X}\boldsymbol{\omega} + \mathbf{m}$
5. Solve Lasso on Page 3, using SGD, ISTA, CD.
6. Plot distance to the optimum objective w.r.t. the CPU time.

Explain how you approximate the optimum.

Consider 4 cases : $(n, p, s) = (200, 200, 0, 5), (200, 200, 0.9), (400, 1500, 0.5), (400, 1500, 0.99)$

Task2.1 : Example

- (Example for one case) Your result may not be the same as mine. Do not worry!



Task 2.2

Regular task 2.2 : Solve Lasso **with** correlation among features

I will only show the part different from task 2.1.

1. Draw $\mathbf{X} \in \mathbb{R}^{n \times p}$, where there are correlations among columns.

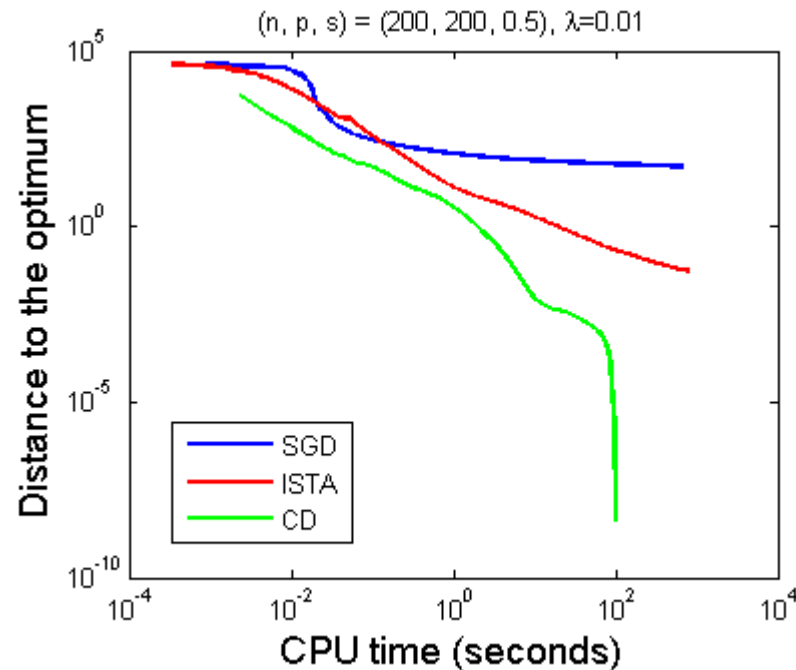
You may do this via the following method.

- (1) Draw $\mathbf{X} \in \mathbb{R}^{n \times p}$, where each element was drawn from $\mathcal{N}(0, \frac{1}{n})$
- (2) Let $\mathbf{C} = p\mathbf{I} + (1 - p)\mathbf{E}$, where $0 \leq p \leq 1$, \mathbf{I} is an identity matrix
 \mathbf{E} is a matrix whose elements are all equal to 1.
- (3) Do Cholesky decomposition : $\mathbf{C} = \mathbf{D}\mathbf{D}^\top$
- (4) let $\mathbf{X} = \mathbf{D}\mathbf{X}$

Consider 4 cases : $(n, p, s) = (200, 200, 0, 5), (200, 200, 0.9),$
 $(400, 1500, 0.5), (400, 1500, 0.99)$

Task2.2 : Example

- (Example for one case) Your result may not be the same as mine. Do not worry!



Extension Task 3

Extension task 3 : Solve Group Lasso with correlation among features

The correlation matrix is block diagonal (block size = group size).

1. Draw $\mathbf{X} \in \mathbb{R}^{n \times p}$, similar as task 2.2 using block diagonal \mathbf{C} .
2. Draw $\boldsymbol{\omega}$ whose sparsity level is s where each non-zero element was drawn from $\mathcal{N}(0, 1)$
3. Draw noise vector \mathbf{m} from i.i.d. Gaussian $\mathcal{N}(0, 0.01 \|\mathbf{X}\boldsymbol{\omega}\|_2^2/n)$
4. Calculate $\mathbf{y} = \mathbf{X}\boldsymbol{\omega} + \mathbf{m}$
5. Solve Group Lasso on Page 8.
6. Plot distance to the optimum objective w.r.t. the CPU time.

Let g_0 be the size for each group. Consider 4 cases

$(n, p, s, g_0) = (200, 200, 0, 5, 10), (200, 200, 0.5, 50), (200, 200, 0.9, 10),$
 $(200, 200, 0.9, 50).$

Tips

- You may consider the following tips when conducting experiments:
 - Tune the hyper-parameter λ to recover the desired result (sparsity pattern in w)
 - Try different strategies in line search (e.g., you may need to tune a and b in the SGD algorithm)
 - Use accurate timers (e.g., 10^{-6} second level accuracy)
 - You may need log-scale to explain the statistics clearly

Bonus

- Here are two problems P1 and P2. Prove that for some T and λ , solving P1 is equivalent to solving P2. You may need to prove the results in two directions.

$$P1 : \min_{\omega} \|\mathbf{y} - \mathbf{X}\omega\|_2^2, \text{ s.t. } \|\omega\|_1 \leq T.$$

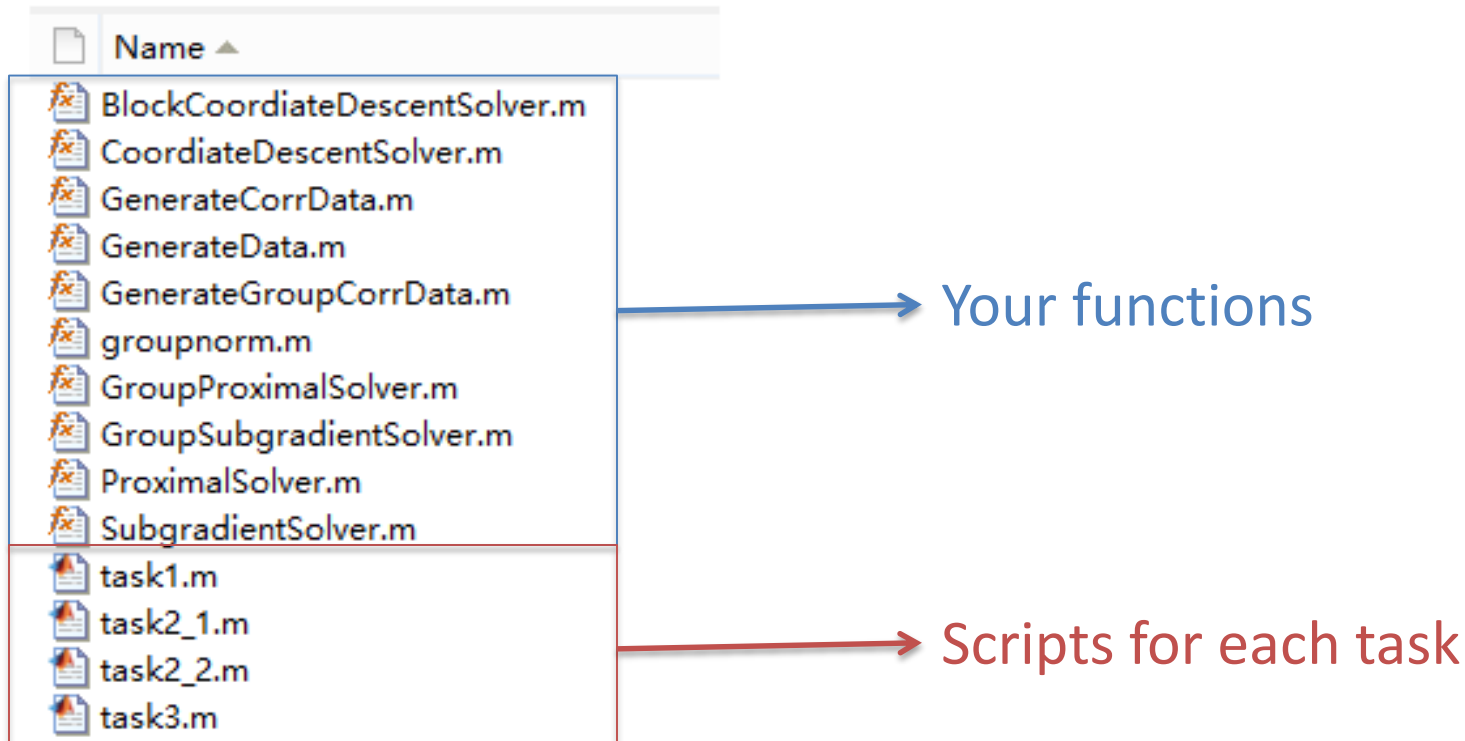
$$P2 : \min_{\omega} \|\mathbf{y} - \mathbf{X}\omega\|_2^2 + \lambda \|\omega\|_1$$

Submission

- Implementation
 - Submit the code implementation before deadline
 - Submit as **.zip/.7z/.tar file**, including:
 - 1) Source Code with necessary comments.
 - 2) Report (.pdf or .doc(x)) containing
 1. Result for each task.
 2. Explain which algorithm is the fastest in your experiments.
 3. If you do the bonus, contain a readable proof in your report.
 - 3) ReadMe explaining
 1. How to run your code for each task (you may use a script).
 2. Personal info (name, class, student id, email).

Code Submission Format

- We recommend you pack your code like this, or you may need to explain how to run your code clearly.



Reference

- F. Bach, R. Jenatton, J. Mairal, G. Obozinski. **Optimization with sparsity-inducing penalties**. *Foundations and Trends in Machine Learning*, 4(1):1-106, 2012.

Thank You!