

Proximal Gradient Descent for the LASSO Problem

Sam Brown Henry Finnila Ryan Voss

June 9, 2025

Contents

1. Introduction

2. Proposed Method

3. Implementation

4. Discussion

5. References

Introduction

Introduction - LASSO Regression

LASSO (Least Absolute Shrinkage and Selection Operator)

- Uses L1 regularization
- Prevents overfitting by shrinking less explanatory coefficients to zero
- Useful for variable selection and model interpretability
- Choose λ through cross-validation
 - Minimize MSE

$$\arg \min_{\beta} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}, \quad \beta \in \mathbb{R}^p \quad (1)$$

LASSO Applications

- Gene selection in high-dimensional microarray and RNA-seq studies [2].
- Feature selection for macroeconomic and financial time series forecasting [3].
- Sparse signal reconstruction in compressed sensing and image processing [4].
- Variable selection in marketing analytics for customer segmentation and churn prediction [5].
- Predictor identification in environmental and climate modeling [6].

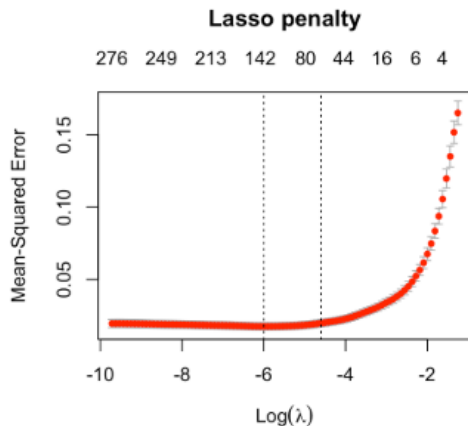


Figure: Minimizing MSE with choice of λ

Ridge vs. LASSO

Ridge Regression:

$$\arg \min_{\beta} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2 \right\} \quad (2)$$

$$\beta = (X^T X + \lambda I)^{-1} X^T y \quad (3)$$

- Can only shrink coefficients asymptotically close to zero
- Better if multicollinearity is present

LASSO Regression:

$$\arg \min_{\beta} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\} \quad (4)$$

- Exact zeros arise from the subgradient interval at zero
- This lets LASSO completely remove certain predictors from a model

Challenges with Optimization

In the LASSO formula, decomposing

$$\arg \min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

into functions

$$g(x) = \frac{1}{2} \|y - X\beta\|_2^2$$

and

$$h(x) = \lambda \|\beta\|_1$$

- L1 penalty, $h(x)$, is not differentiable
- Standard gradient descent fails
- Need a new method...

Proposed Method

Proximal/Generalized Gradient Descent

Consider $f : \mathbb{R}^p \rightarrow \mathbb{R}$ that is not necessarily differentiable, but can be decomposed

$$f(x) = g(x) + h(x)$$

- g convex and differentiable
- h convex

LASSO objective can easily be decomposed in this manner.

$$\underbrace{\frac{1}{2}\|y - X\beta\|_2^2}_{g(x)} + \underbrace{\lambda\|\beta\|_1}_{h(x)}$$

Derivation (1/2)

Recall gradient descent update with step size α for differentiable g :

$$x_{k+1} = x_k - \alpha \nabla g(x_k)$$

Comes from quadratic approximation

$$x_{k+1} = \arg \min_z \left(g(x_k) + \nabla g(x_k)^\top (z - x_k) + \frac{1}{2\alpha} \|z - x_k\|_2^2 \right)$$

Applying this approximation to g gives PGD update

$$x_{k+1} = \arg \min_z \left(g(x_k) + \nabla g(x_k)^\top (z - x_k) + \frac{1}{2\alpha} \|z - x_k\|_2^2 + h(z) \right)$$

Derivation (2/2)

Equivalently,

$$x_{k+1} = \arg \min_z \left(\frac{1}{2\alpha} \|z - (x_k - \alpha \nabla g(x_k))\|_2^2 + h(z) \right)$$

- Remain close to GD update for g
- Enforce small h

The Proximal Operator

Define $\text{prox}_{\alpha,h} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ by

$$\text{prox}_{\alpha,h}(x) = \arg \min_{z \in \mathbb{R}^p} \left(\frac{1}{2\alpha} \|z - x\|_2^2 + h(z) \right).$$

Note that proximal mapping depends on α and h but not g . Our update,

$$x_{k+1} = \arg \min_z \left(\frac{1}{2\alpha} \|z - (x_k - \alpha \nabla g(x_k))\|_2^2 + h(z) \right),$$

becomes

$$x_{k+1} = \text{prox}_{\alpha,h}(x_k - \alpha \nabla g(x_k)).$$

How is This Useful?

$$x_{k+1} = \text{prox}_{\alpha,h}(x_k - \alpha \nabla g(x_k))$$
$$\text{prox}_{\alpha,h}(x) = \arg \min_z \left(\frac{1}{2\alpha} \|z - x\|_2^2 + h(z) \right)$$

- Goal: Solve a minimization problem
- Proposed method: Solve a minimization problem on every iteration?
- Need an easy way to evaluate $\text{prox}_{\alpha,h}$
- Possible for many important h functions, including L_1 regularization (LASSO objective)

Proximal Mapping for the LASSO Problem

For the LASSO objective,

$$h(\beta) = \lambda \|\beta\|_1.$$

In this case,

$$\text{prox}_{\alpha, h}(\beta) = S_{\alpha\lambda}(\beta),$$

where $S_{\alpha\lambda} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is the soft thresholding operator defined element wise by

$$[S_{\alpha\lambda}(\beta)]_i = \begin{cases} \beta_i + \alpha\lambda & \text{if } \beta_i < -\alpha\lambda \\ 0 & \text{if } -\alpha\lambda < \beta_i < \alpha\lambda \\ \beta_i - \alpha\lambda & \text{if } \alpha\lambda < \beta_i, \end{cases}$$

for $i = 1, \dots, p$. This allows us to use proximal gradient descent for the LASSO problem. The resulting algorithm is often called the Iterative Soft Thresholding Algorithm (ISTA).

Convergence Guarantees

Theorem

Suppose that

- $f : \mathbb{R}^p \rightarrow \mathbb{R}$ can be decomposed into $f(x) = h(x) + g(x)$
- g is convex, differentiable, and L -smooth
- h is convex
- $\text{prox}_{\alpha, h}$ can be evaluated.

Then, proximal gradient descent with step size $0 < \alpha \leq \frac{1}{L}$ satisfies

$$|f(x_k) - f(x^*)| \leq \frac{\|x_0 - x^*\|_2^2}{2\alpha k}.$$

Remarks:

- Convergence is $\mathcal{O}(1/k)$
- Same convergence rate as gradient descent, but we are just counting iterations

Implementation

Implementation

To demonstrate the implementation of proximal gradient descent, we will solve the well-known LASSO (Least Absolute Shrinkage and Selection Operator) problem. As previously discussed, LASSO requires solving the following optimization problem:

$$\arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

This problem is not differentiable due to the ℓ_1 -norm term $\lambda \|\beta\|_1$, and therefore standard gradient descent cannot be directly applied. Instead, we use **proximal gradient descent**, which is specifically designed to handle non-smooth regularization terms.

We implement this optimization algorithm in Python using the PyTorch library.

Data Initialization

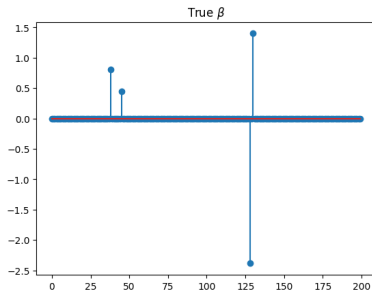
For this implementation, we use synthetic data to keep things simple and easy to control. The code below creates a high-dimensional dataset with a sparse true coefficient vector `betaTrue`, which works well for testing Lasso regression.

Listing 1: Synthetic data creation

```
1 import torch
2 import matplotlib.pyplot as plt
3 torch.manual_seed(1)
4 # synthetic data
5 N = 40      # training examples
6 d = 200     # number of features
7 nnz = 4     # number of nonzero components
8 m = torch.randperm(d) # random permutation of indices
9 betaTrue = torch.zeros(d) # initialize beta vector
10 betaTrue[m[:nnz]] = 5 * torch.randn(nnz) # assign 4 random values
11 X = torch.randn(N, d) # feature matrix
12 ones = torch.ones(X.size(0), 1) # for intercept (if necessary)
13 X_new = torch.cat((ones, X), dim=1) # combine intercept column with X
14 noise = 0.1 * torch.randn(N)
15 y = X @ betaTrue + noise
```

Motivation for Using Lasso

- Lasso fits our setup with only a few nonzero coefficients.
- nnz sets how many predictors matter; noise is added to y .
- Plotting β_{True} shows the sparsity Lasso aims to recover.



Defining the Proximal Operator

To solve the Lasso problem, we must handle the non-differentiable ℓ_1 penalty term. This is done using the **proximal operator** for the L1 norm, which shrinks each component of a vector toward zero by a fixed amount, promoting sparsity.

We implement this using a `soft_threshold` function, which moves each value closer to zero by `alpha`, and sets small values exactly to zero. The `proxL1Norm` function applies this elementwise and optionally skips the first entry (commonly used for intercept terms).

The operator corresponds to the soft-thresholding function:

$$[S_{\alpha\lambda}(\beta)]_i = \begin{cases} \beta_i + \alpha\lambda & \text{if } \beta_i < -\alpha\lambda \\ 0 & \text{if } -\alpha\lambda < \beta_i < \alpha\lambda \\ \beta_i - \alpha\lambda & \text{if } \beta_i > \alpha\lambda \end{cases}$$

Proximal Operator Implementation

The code below implements the soft-thresholding operator using PyTorch's `clamp` function to enforce sparsity in Lasso regression.

```
1 def soft_threshold(beta, alpha):
2     # shrinks each value toward zero
3     return torch.sign(beta) * torch.clamp(torch.abs(beta) - alpha, min=0.0)
4
5 def proxL1Norm(betaHat, alpha, penalizeALL=True):
6     # applies soft-thresholding; optionally skips the intercept
7     out = soft_threshold(betaHat, alpha)
8     if not penalizeALL:
9         out[0] = betaHat[0]
10    return out
```

This operation matches the mathematical definition of $S_{\alpha\lambda}(\beta)$ and is key to enabling Lasso's sparse recovery.

Proximal Gradient Descent Implementation

```
1 def solveLasso_proxGrad(X, y, lambda):
2     maxIter = 300
3     alpha = 0.005 # step size
4
5     N, d = X.shape # number of samples, number of features
6     beta = torch.zeros(d, dtype=torch.float32) # initialize
7     costFunVals = torch.zeros(maxIter) # initialize
8     for t in range(maxIter):
9         grad = X.T @ (X @ beta - y) # compute gradient of differentiable term
10        beta = proxL1Norm(beta - alpha * grad, alpha * lambda) # gradient step
11                               then apply prox
12
13        residual = X @ beta - y # compute and store residual
14        cost = 0.5 * torch.norm(residual)**2 + lambda * torch.sum(torch.abs(beta)
15                               )
16        costFunVals[t] = cost
17
18        print(f"Iteration: {t}, Objective function value: {cost.item():.4f}")
19    return beta, costFunVals
```

Proximal Gradient Update Step

The core of the `solveLasso_proxGrad` function is the following line:

```
1 beta = proxL1Norm(beta - alpha * grad, alpha * lambda)
```

This line performs the proximal gradient update:

- First, it takes a gradient step on the smooth squared loss.
- Then, it applies the proximal operator to handle the non-smooth ℓ_1 penalty.

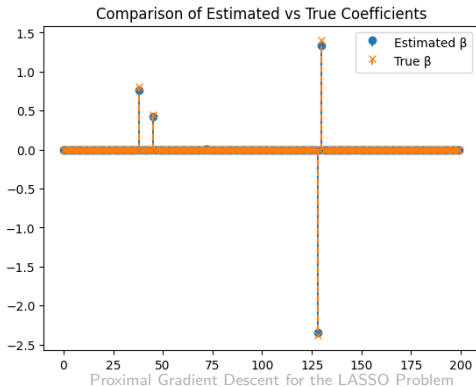
$$\beta^{(k+1)} = \text{prox}_{\alpha, h} \left(\beta^{(k)} - \alpha_k \nabla f \left(\beta^{(k)} \right) \right)$$

The second argument in `proxL1Norm` is the scaled regularization parameter $\alpha \cdot \lambda$. As shown earlier, this causes each component of β to be shifted toward zero by this exact amount, as defined by the soft-thresholding operator.

Testing the Implementation

To test our proximal gradient descent code, we call the function with our feature matrix, target vector, and a specified λ value:

```
1 lambda = 2
2 beta, costFunVals = solveLasso_proxGrad(X, y, lambda)
```



Lasso with scikit-learn (1/2)

- We now use the built-in Lasso class from scikit-learn.
- This is more efficient and comes with tools like preprocessing and cross-validation.
- We use the `california_housing` dataset, which is convenient and well-supported.

Import libraries and prepare data:

```
1 from sklearn.datasets import fetch_california_housing
2 from sklearn.linear_model import Lasso
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5
6 # Load and split data
7 X, y = fetch_california_housing(return_X_y=True)
8 X_train, X_test, y_train, y_test = train_test_split(X, y)
9
10 # Scale features
11 scaler = StandardScaler()
12 X_train = scaler.fit_transform(X_train)
13 X_test = scaler.transform(X_test)
```

Lasso with scikit-learn (2/2)

- We initialize the Lasso model and train it on the scaled training data.
- Then, we evaluate it on the test set using Mean Squared Error (MSE).

Model training and evaluation:

```
1 # Initialize and train Lasso model
2 lasso = Lasso(alpha=1.0)
3 lasso.fit(X_train, y_train)
4
5 # Predict and evaluate
6 y_pred = lasso.predict(X_test)
7 from sklearn.metrics import mean_squared_error
8 mse = mean_squared_error(y_test, y_pred)
9 print(f"Mean Squared Error: {mse}")
```

The MSE gives a quantitative measure of the model's prediction error on unseen data.

Discussion

Discussion of our Algorithm

Strengths

- **Sparsity promotion.** The soft-thresholding operator can set coefficients to zero
- **Simplicity and ease of implementation.**
- **Guaranteed convergence.** Under standard assumptions ISTA converges at a $\mathcal{O}(1/k)$ rate .

Weaknesses

- **Slow convergence.** The $\mathcal{O}(1/k)$ rate can be slow in practice
- **Step-size sensitivity.** Tuning the parameters it one of the most crucial, yet difficult parts of the process.

Thank you!

References

References



R. Tibshirani, *Convex Optimization: Proximal Gradient Methods*, Lecture 8, Carnegie Mellon Univ., Pittsburgh, PA, USA, 2015. [Online]. Available: <https://www.stat.cmu.edu/~ryantibs/convexopt-S15/lectures/08-prox-grad.pdf>



F. Hu, R. Zhao, and J. Xu, "Supervised group Lasso with applications to microarray data analysis," *BMC Bioinformatics*, vol. 9, p. 10, 2008.



A. Bonavito and L. Morelli, "Forecasting macroeconomic time series: Lasso-based approaches," *Econometric Reviews*, vol. 33, no. 5–6, pp. 594–616, 2014.



D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.



L. Wang, H. Li, and M. Zhang, "Customer churn prediction based on Lasso and random forest models," in *Proc. 2018 IEEE Int. Conf. Cloud Comput. Big Data Anal.*, Chengdu, China, 2018, pp. 231–235.



P. Liang, K. Fraedrich, and R. Schnur, "Predictor selection for downscaling GCM data with Lasso," *J. Geophys. Res. Atmos.*, vol. 117, no. D11, 2012.