

High-Dimensional Regression: LASSO and AdaBoost

Abdallah Abdelsameia, Louis Rodriguez

Toulouse School of Economics - M2 Data Science for Social Sciences

March 2025

1 Introduction

2 LASSO Implementations

- Coordinate Descent (CD)
- ISTA
- FISTA
- Square-root LASSO

3 Experimental Scenarios

- Scenario 1: Baseline
- Scenario 2: Toeplitz Correlation
- Scenario 3: High-Dimensional Scale-Up

4 Adaboost Implementations

- MyAdaBoost (Manual Implementation)
- AdaBoost: Scikit-Learn Implementation

5 Comparison and Conclusion

- **High-dimensional regression:** We often have $p \gg n$, or at least p comparable to n , where many features may be redundant.
- **LASSO** addresses this by:
 - Minimizing a penalized least squares objective:

$$\min_{b \in \mathbb{R}^p} \frac{1}{n} \|Y - Xb\|_2^2 + 2\tau \|b\|_1,$$

where $\|b\|_1 = \sum_j |b_j|$.

- Encouraging *sparsity* in b : many coefficients become zero.
 - Simultaneously performing variable selection and shrinkage.
- In this presentation, we focus on **four LASSO variants**:
 - Coordinate Descent (CD)
 - ISTA (Iterative Shrinkage-Thresholding)
 - FISTA (Fast ISTA with momentum)
 - Square-root Lasso (scaled Lasso)
- **We run experiments** on synthetic data to see:
 - How these methods behave under different design conditions (correlation, dimensionality),
 - Their ability to recover the true sparse support,
 - Their runtime and convergence properties.

LASSO: Coordinate Descent (Method 1)

```
1 import numpy as np
2
3 def soft_thresh(r, lam):
4     return np.sign(r) * max(abs(r) - lam, 0)
5
6 def lasso_cd(X, y, lam, max_iter=5000, tol=1e-4):
7     n, p = X.shape
8     b = np.zeros(p)
9     for it in range(max_iter):
10         b_old = b.copy()
11         for j in range(p):
12             r = y - X @ b + b[j] * X[:, j]
13             rho = (X[:, j] @ r) / n
14             norm_j = (X[:, j]**2).sum() / n
15             b[j] = soft_thresh(rho, lam) / norm_j
16         if np.linalg.norm(b - b_old) < tol:
17             break
18     return b
```

Implementation Description:

- Initialize β to zeros.
- **Loop** over a maximum of `max_iter`:
 - For each coordinate j :
 - 1 Compute a *partial residual*
 $r = y - X\beta + \beta_j X_{:,j}$.
 - 2 Find correlation ρ between r and $X_{:,j}$.
 - 3 Compute $b[j]$ by *soft thresholding* ρ , scaled by the feature norm.
 - Check if the updated β changes less than `tol`.
- Return β when converged or after max iterations.

```

1 def soft_thresh(r, lam):
2     return np.sign(r) * max(abs(r) - lam, 0)
3
4 def lasso_ista(X, y, lam, L=None,
5               max_iter=1e4, tol=1e-4):
6     n, p = X.shape
7     if L is None:
8         L = (np.linalg.norm(X, 2)**2)/n
9     b = np.zeros(p)
10    for k in range(int(max_iter)):
11        b_old = b.copy()
12        grad = (X.T @ (X @ b - y)) / n
13        b = soft_thresh(b - grad / L, lam / L)
14        if np.linalg.norm(b - b_old) < tol:
15            break
16    return b

```

Implementation Description:

- Initialize β as zeros.
- Estimate Lipschitz constant $L \approx \frac{\|X\|_2^2}{n}$ if not provided.
- For each iteration:
 - 1 Compute gradient of the least-squares term: $\nabla f(\beta)$.
 - 2 Perform *proximal* update:

$$\beta \leftarrow \text{soft_thresh}\left(\beta - \frac{1}{L} \nabla f, \frac{\lambda}{L}\right).$$
 - 3 Check convergence on $\|\beta - \beta_{\text{old}}\|$.
- Return β .

```
1 def soft_thresh(r, lam):
2     return np.sign(r) * max(abs(r) - lam, 0)
3
4 def lasso_fista(X, y, lam, L=None,
5               max_iter=1e4, tol=1e-4):
6     n, p = X.shape
7     if L is None:
8         L = (np.linalg.norm(X, 2)**2)/n
9     b = np.zeros(p)
10    yk = b.copy()
11    t = 1.0
12    for k in range(int(max_iter)):
13        b_old = b.copy()
14        grad = (X.T @ (X @ yk - y)) / n
15        b = soft_thresh(yk - grad / L, lam / L)
16        t_new = 0.5 * (1 + (1 + 4*t*t)**0.5)
17        yk = b + ((t - 1)/t_new) * (b - b_old)
18        t = t_new
19        if np.linalg.norm(b - b_old) < tol:
20            break
21    return b
```

Implementation Description:

- Similar to ISTA but uses an **acceleration** (Nesterov momentum).
- We keep track of:
 - A separate y_k vector (the “extrapolated” point),
 - A momentum scalar t updated each iteration.
- Each iteration:
 - 1 Gradient step on y_k .
 - 2 Apply soft-threshold with λ/L .
 - 3 Update t_{k+1} and combine the new b with old b to form y_k .
- Typically faster convergence than ISTA.

Square-root Lasso (Method 4)

```
1 def sqrt_lasso(X, y, tau,
2               max_outer=100, tol=1e-4):
3     n, p = X.shape
4     sigma = np.linalg.norm(y)/np.sqrt(n)
5     b = np.zeros(p)
6     for it in range(max_outer):
7         b_old, s_old = b.copy(), sigma
8         # solve LASSO subproblem
9         b = lasso_fista(X, y, lam=tau*sigma)
10        r = y - X@b
11        sigma = np.linalg.norm(r)/np.sqrt(n)
12        if (abs(sigma - s_old)<tol and
13            np.linalg.norm(b-b_old)<tol):
14            break
15    return b, sigma
```

Implementation Description:

- We minimize:

$$\min_{b, \sigma > 0} \frac{\sigma}{2} + \frac{\|y - Xb\|^2}{2\sigma n} + \tau \|b\|_1.$$

- Initialize σ from $\|y\|/\sqrt{n}$.
- Outer loop:
 - 1 Fix σ , solve standard LASSO with $\lambda = \tau \cdot \sigma$ (using FISTA).
 - 2 Update σ from new residual $\|y - Xb\|/\sqrt{n}$.
 - 3 Check if both b and σ changed below tol .
- Returns (b, σ) upon convergence.

- We design various synthetic setups to test LASSO performance:
 - ① **Baseline (make_regression):** Low correlation, moderate dimension.
 - ② **Toeplitz Correlated:** Vary correlation ρ .
 - ③ **High-Dimensional Scale-Up:** Fix n , grow p .
- Compare methods on:
 - Support recovery
 - False positives
 - Runtime
- Next slides: results and highlights for each scenario.

make_regression: Key Highlights

Generate Synthetic Regression Data

What It Does:

- Generates data using a random linear model.
- Constructs outputs from a subset of `n_informative` features.
- Optionally adds Gaussian noise to simulate real-world data.
- Supports both well-conditioned and low-rank data (via `effective_rank`).

Important Parameters:

- `n_samples`: Number of samples.
- `n_features`: Total features.
- `n_informative`: Features used to generate the output.
- `noise`: The standard deviation of the gaussian noise applied to the output.

Example Usage:

```
1 from sklearn.datasets import make_regression
2
3 # Generate data with 100 samples, 100 features, 10
  informative, and noise level 10
4 X, y = make_regression(n_samples=100, n_features=100,
5                       n_informative=10, noise=10,
6                       random_state=42)
```

Scenario 1: Baseline (`make_regression`)

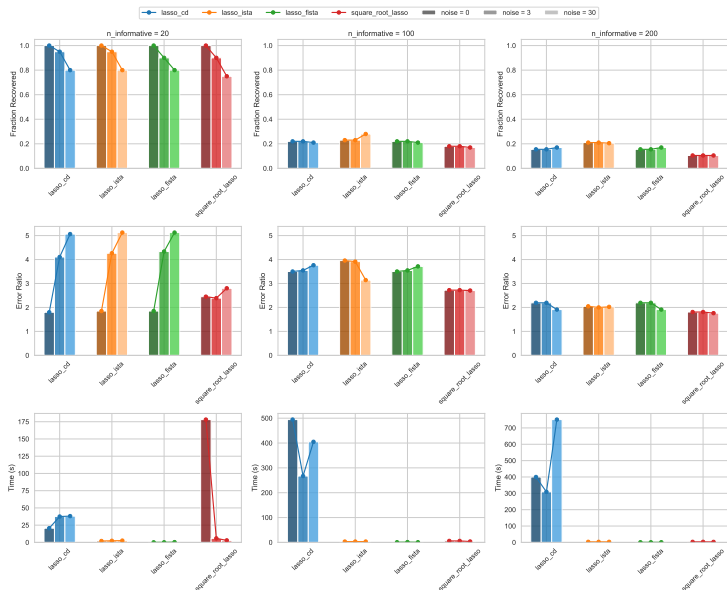
- **Data Generation:**

- Use `sklearn.datasets.make_regression` with:
 - $n = 100$ samples, $p = 1000$ features (example),
 - `n_informative` values in $\{20, 100, 200\}$,
 - noise values in $\{0, 3, 30\}$,
 - Sparse coefficients.
- Standardize X and center y .

- **Expectations:**

- Lasso can achieve near-oracle performance if the design is well-behaved (i.e., no extreme correlations) and the true β is sparse.
- `make_regression` by default produces a "nice" design, typically satisfying the restricted eigenvalue condition.
- As noise increases, Lasso's support recovery (identifying nonzero coefficients) becomes more challenging.
- Error bounds include a penalty factor that grows with the noise variance.

Scenario 1: Results



Scenario 2: Correlated Design (Toeplitz)

- **Data Generation:**

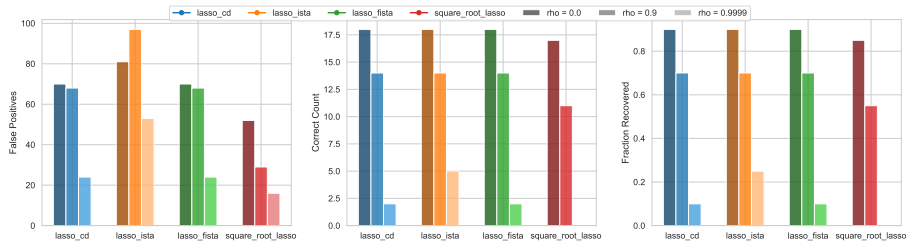
- Covariance Σ with $\Sigma_{ij} = \rho^{|i-j|}$.
- Sample $X_i \sim \mathcal{N}(0, \Sigma)$, $i = 1, \dots, n$.
- True β sparse, $y = X\beta + \epsilon$.

- **Vary** $\rho \in \{0.0, 0.9, 0.9999\}$ to see effect of correlation.

- **Expectations:**

- Lasso can achieve near-oracle performance if the design is well-behaved (i.e., no extreme correlations) and the true β is sparse.
- When correlation increases the RE (or RIP) assumptions can fail, which often hurts Lasso's ability to correctly recover the sparse support.
- We expect to see more false positives because different columns in the same "correlated cluster" can all partially explain the same signal.
- Convergence might slow down for iterative methods if the design is "ill-conditioned." and we expect to see longer running times.

Scenario 2: Results



Parameters: $n = 100$, $p = 1000$, $\rho \in \{0.0, 0.9, 0.9999\}$, $n_{\text{informative}} = 20$, $\sigma_{\text{noise}} = 1.0$, $\lambda = 0.1$, $\tau = 0.2$,
 $\text{random_state}=42$.

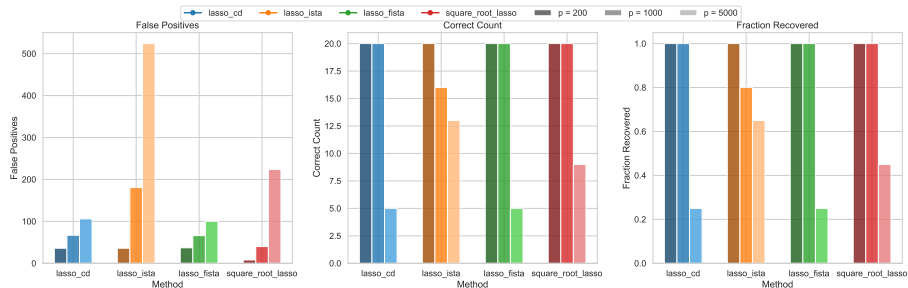
- **Setup:**

- Fix $n = 100$, let $p \in \{200, 1000, 5000\}$ or more.
- Keep $n_{\text{informative}} = 20$ and moderate noise.
- Evaluate each method's runtime, fraction recovered, etc.

- **Expectations:**

- As p grows, naive coordinate descent might slow.
- $\log(p)$ factor in theory: can see more false positives if λ not tuned carefully.
- It could be harder for lasso to distinguish the true signal source if some of the features are correlated by chance.

Scenario 3: Results



Parameters: $p_{\text{list}} = [200, 1000, 5000]$, $n = 100$, $n_{\text{informative}} = 20$, $\sigma_{\text{noise}} = 1.0$, $\lambda = 0.1$, $\tau = 0.2$,
 $\text{random_state}=42$.

- LASSO solutions strongly depend on design matrix properties:
 - Correlation, rank deficiencies, high p all degrade success of feature recovery.
- FISTA often converges faster than ISTA; coordinate descent can be effective in low correlation.
- Square-root Lasso can help handle unknown noise level automatically.

- AdaBoost: adaptive boosting for feature selection and regression.
- 2 versions of AdaBoost for regression:
 - Scikit-learn AdaBoostRegressor.
 - Custom MyAdaBoostRegressor implementation.
- Experiments on synthetic data follow.

```
1 def fit(self, X, y):
2     n_samples = X.shape[0]
3     sample_weights = np.ones(n_samples) / n_samples
4
5     for m in range(self.n_estimators):
6         estimator = deepcopy(DecisionTreeRegressor(max_depth=
7             self.max_depth)) # Stronger weak learners
8         estimator.fit(X, y, sample_weight=sample_weights)
9         y_pred = estimator.predict(X)
10
11        # Compute absolute errors and median threshold
12        errors = np.abs(y - y_pred)
13        tau = np.median(errors)
14        L = (errors > tau).astype(int) # Binary loss
15
16        # Compute weighted error
17        err_m = np.sum(sample_weights * L) / np.sum(
18            sample_weights)
19        err_m = np.clip(err_m, 1e-10, 0.4999) # Prevent
20            division instability
21
22        # Compute model weight (alpha)
23        alpha_m = 0.5 * np.log((1 - err_m) / err_m) # Scaling
24            alpha
25        beta_m = err_m / (1 - err_m) # Compute beta_m
26
27        # Update sample weights using beta_m^(1 - L)
28        sample_weights **= beta_m ** (1 - L)
29        sample_weights /= np.sum(sample_weights)
30        self.estimators_.append(estimator)
31        self.estimator_weights_.append(alpha_m)
32
33    def predict(self, X):
34        predictions = np.array([
35            alpha * est.predict(X) for alpha, est in zip(self.
36                estimator_weights_, self.estimators_)
37        ])
38    return np.sum(predictions, axis=0) # Weighted sum
```

Key Points:

- Iteratively trains multiple DecisionTreeRegressors and assigns higher weights to harder-to-predict samples.
- Weak learners are weighted based on their error rate; better models get higher weights.
- The error err_m is computed as a weighted error rate.
- The model weight α_m is calculated as:

$$\alpha_m = \frac{1}{2} \log \left(\frac{1 - err_m}{err_m} \right)$$

- The beta value (β_m) used for weight updates is:

$$\beta_m = \frac{err_m}{1 - err_m}$$

- The weight update rule is:

$$w_i = w_i \times \beta_m^{(1-L)}$$

```
1 from sklearn.ensemble import AdaBoostRegressor
2 from sklearn.tree import DecisionTreeRegressor
3
4 # Create AdaBoost with decision stumps
5 ada_reg = AdaBoostRegressor(
6     estimator=DecisionTreeRegressor(max_depth=4),
7     n_estimators=50,
8     learning_rate=1.0,
9     random_state=42
10 )
11
12 # Train the model
13 ada_reg.fit(X_train, y_train)
14
15 # Predictions
16 y_pred = ada_reg.predict(X_test)
```

Key Points:

- The weak learner used is a `DecisionTreeRegressor` with `max_depth=4`, meaning each tree can have a maximum depth of 4.
- Unlike classification AdaBoost, the regression version uses weighted residuals to improve performance.
- The model is trained sequentially, where each new tree focuses on correcting the errors of the previous trees.
- `n_estimators=50`: Uses 50 weak learners (decision trees).
- `learning_rate=1.0`: Controls how much influence each tree has on the final prediction.

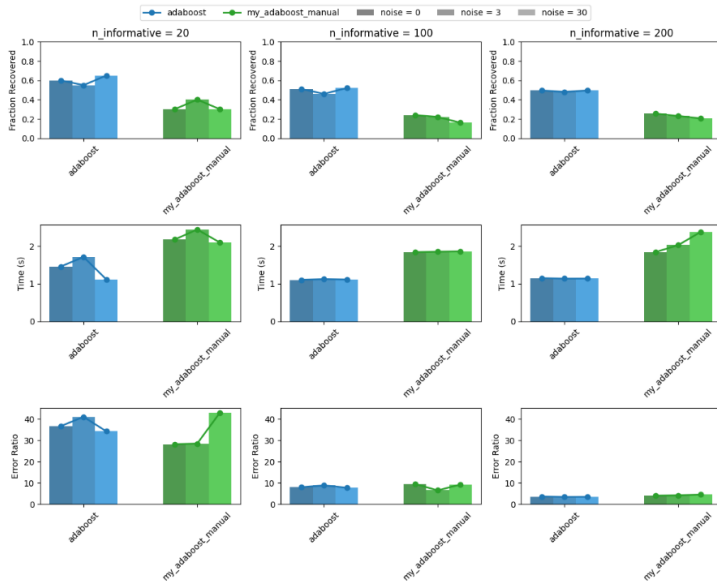
- **Data Generation:**

- Use `sklearn.datasets.make_regression` with:
 - $n = 100$ samples, $p = 1000$ features.
 - `n_informative = 20, 100, 200`.
 - noise values in $\{0, 3, 30\}$.
 - Sparse coefficients.
- Standardize X and center y .
- Split into training (80%) and test (20%).

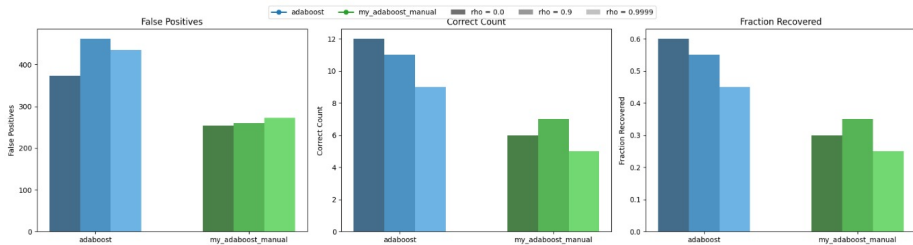
- **Expectations:**

- AdaBoost assigns more weight to samples that are harder to predict.
- Feature importance estimation depends on boosting weight updates.
- Higher noise levels make it difficult for weak learners to detect patterns.
- Manual implementation should perform similarly to scikit-learn's, with minor differences due to weight updates.

Scenario 1 on AdaBoost



Scenario 2: Adaboost - Correlated Design (Toeplitz)



Comparison Table

Aspect	Lasso Regression	Adaboost Regression
Type	Linear model	Ensemble method
Regularization	L1 (sparsity)	No direct regularization
Strength	Feature selection	Captures complex non-linearity
Weakness	Struggles with non-linearity	Sensitive to outliers
Use Case	High-dimensional, sparse data	Complex relationships, boosting weak models

[Table:](#) Comparison of Lasso and Adaboost Regression



Belloni, A., Chernozhukov, V., & Wang, L. (2011).

"Square-root Lasso: Pivotal recovery of sparse signals via conic programming."

Biometrika, 98(4), 791–806.

(**Square-root Lasso**)



Beck, A., & Teboulle, M. (2009).

"A fast iterative shrinkage-thresholding algorithm for linear inverse problems."

SIAM Journal on Imaging Sciences, 2(1), 183–202.

(**FISTA Algorithm Paper**)



Friedman, J., Hastie, T., & Tibshirani, R. (2010).

"Regularization Paths for Generalized Linear Models via Coordinate Descent."

Journal of Statistical Software, 33(1), 1–22.

(**Coordinate Descent for Lasso**)



Christophe Giraud. *Introduction to High-Dimensional Statistics*. CRC Press, 2021, Chapter 11. <https://www.imo.universite-paris-saclay.fr/~christophe.giraud/Orsay/Bookv3.pdf>



Our Implementations:

Coordinate Descent, ISTA, FISTA, and Square-root Lasso in Python:

<https://github.com/aabdelsameia1/HighDimensionalProject>