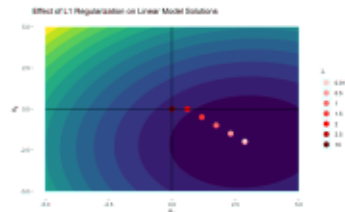


Introduction to Machine Learning

Regularization

Lasso Regression

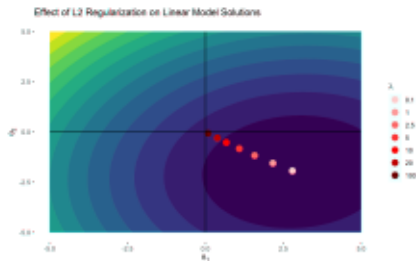
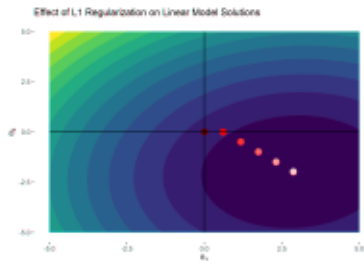


Learning goals

- Lasso regression / L_1 penalty
- Know that lasso selects features
- Support recovery

LASSO REGRESSION / 2

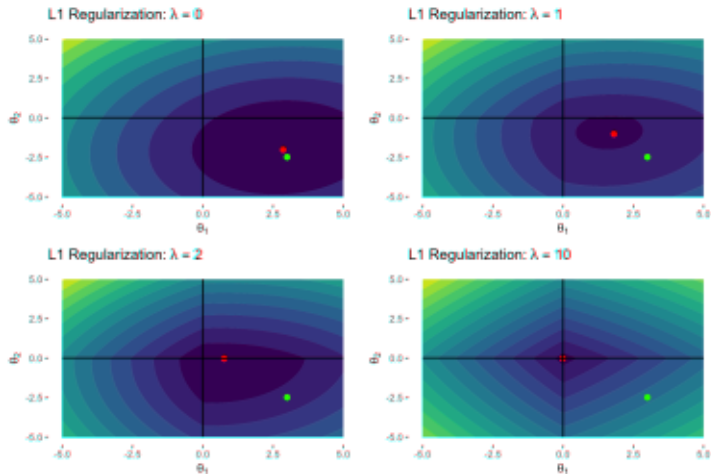
Let $y = 3x_1 - 2x_2 + \epsilon$, $\epsilon \sim N(0, 1)$. The true minimizer is $\theta^* = (3, -2)^T$. LHS = L1 regularization; RHS = L2



With increasing regularization, $\hat{\theta}_{lasso}$ is pulled back to the origin, but takes a different "route". θ_2 eventually becomes 0!

LASSO REGRESSION / 3

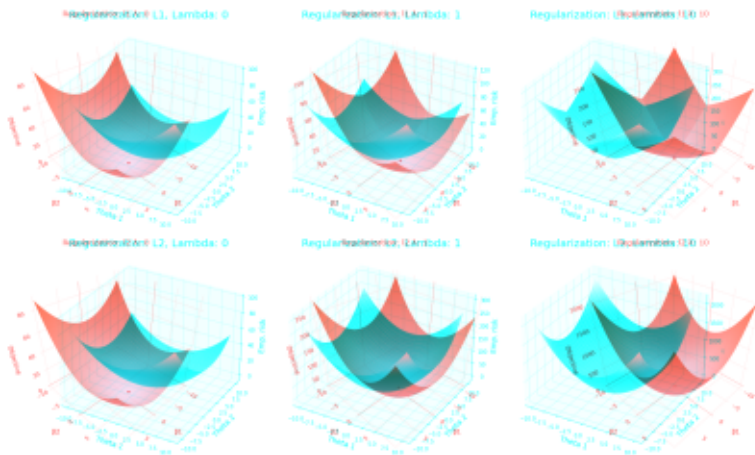
Contours of regularized objective for different λ values.



Green = true minimizer of the unreg. objective and red = lasso solution.

LASSO REGRESSION / 4

Regularized empirical risk $\mathcal{R}_{\text{reg}}(\theta_1, \theta_2)$ using squared loss for $\lambda \uparrow$. L_1 penalty makes non-smooth kinks at coordinate axes more pronounced, while L_2 penalty warps \mathcal{R}_{reg} toward a "basin" (elliptic paraboloid).

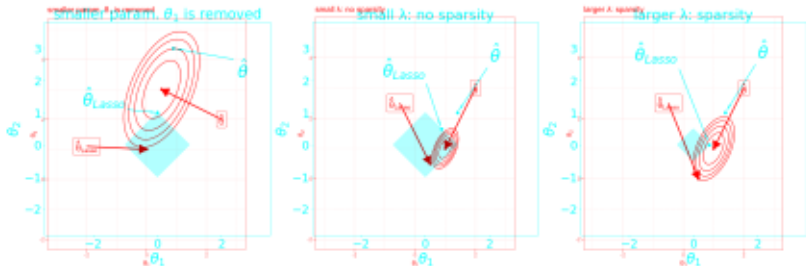


LASSO REGRESSION / 5

We can also rewrite this as a constrained optimization problem. The penalty results in the constrained region to look like a diamond shape.

$$\min_{\theta} \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)} | \theta) \right)^2 \text{ subject to: } \|\theta\|_1 \leq t$$

The kinks in L_1 enforce sparse solutions because “the loss contours first hit the sharp corners of the constraint” at coordinate axes where **(some) entries are zero**.



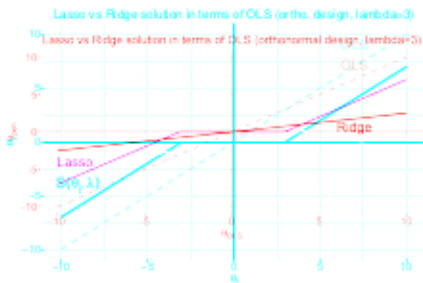
L1 AND L2 REG. WITH ORTHONORMAL DESIGN

For special case of orthonormal design $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ we can derive a closed-form solution in terms of $\hat{\theta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{y}$:

$$\hat{\theta}_{Lasso} = \text{sign}(\hat{\theta}_{OLS}) (|\hat{\theta}_{OLS}| - \lambda)_+ \quad (\text{sparsity})$$

Function $S(\theta; \lambda) := \text{sign}(\theta) (|\theta| - \lambda)_+$ is called **soft thresholding operator**:
For $|\theta| \leq \lambda$ it returns 0, whereas params $|\theta| > \lambda$ are shrunk toward 0 by λ .
Comparing this to $\hat{\theta}_{Ridge}$ under orthonormal design:

$$\hat{\theta}_{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = ((1 + \lambda) \mathbf{I})^{-1} \hat{\theta}_{OLS} = \frac{\hat{\theta}_{OLS}}{1 + \lambda} \quad (\text{no sparsity})$$



COMPARING SOLUTION PATHS FOR L1/L2

- Ridge results in smooth solution path with non-sparse params
- Lasso induces sparsity, but only for large enough λ

