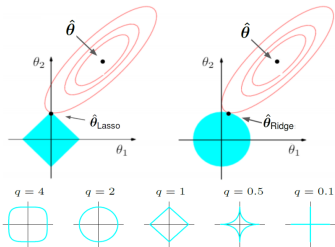
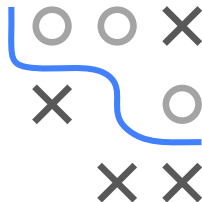


Introduction to Machine Learning

Regularization

Other Regularizers

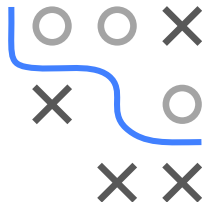
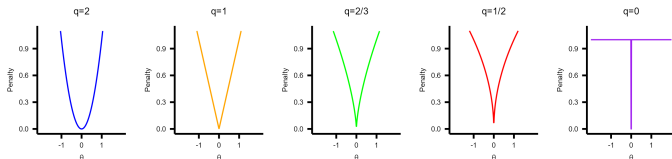


Learning goals

- L1/L2 regularization induces bias
- L_q (quasi-)norm regularization
- L0 regularization
- SCAD and MCP

L0 REGULARIZATION

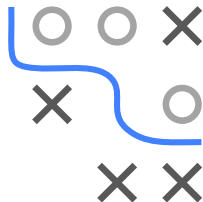
$$\mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) = \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_0 := \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \sum_j |\theta_j|^0.$$



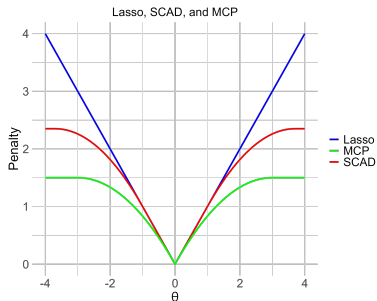
- L0 "norm" simply counts the nr of non-zero params
- Induces sparsity more aggressively than L_1 , but does not shrink
- AIC and BIC are special cases of L_0
- L_0 -regularized risk is not continuous or convex
- NP-hard to optimize; for smaller n and p somewhat tractable, efficient approximations are still current research

Smoothly Clipped Absolute Deviations:
non-convex, $\gamma > 2$ controls how fast penalty “tapers off”

$$\text{SCAD}(\theta \mid \lambda, \gamma) = \begin{cases} \lambda|\theta| & \text{if } |\theta| \leq \lambda \\ \frac{2\gamma\lambda|\theta| - \theta^2 - \lambda^2}{2(\gamma-1)} & \text{if } \lambda < |\theta| < \gamma\lambda \\ \frac{\lambda^2(\gamma+1)}{2} & \text{if } |\theta| \geq \gamma\lambda \end{cases}$$



- Lasso, quadratic, then const
- Smooth
- Contrary to lasso/ridge, SCAD continuously relaxes penalization rate as $|\theta|$ increases above λ

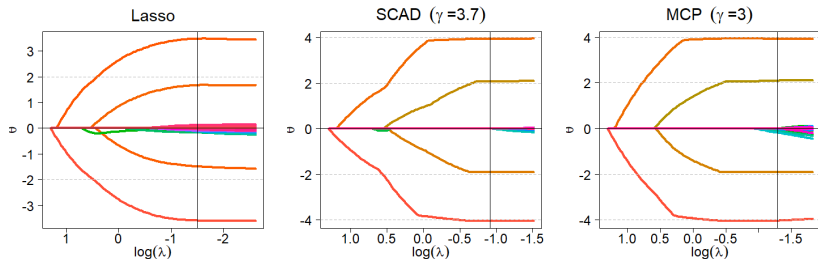


EXAMPLE: COMPARING REGULARIZERS

Let's compare coeff paths for lasso, SCAD, and MCP.

We simulate $n = 100$ samples from the following DGP:

$$\mathbf{y} = \mathbf{x}^\top \boldsymbol{\theta} + \varepsilon, \quad \boldsymbol{\theta} = (4, -4, -2, 2, 0, \dots, 0)^\top \in \mathbb{R}^{1500}, \quad x_j, \varepsilon \sim \mathcal{N}(0, 1)$$



Vertical lines mark optimal λ from 10CV.

Conclusion: Lasso underestimates true coeffs while SCAD/MCP achieve unbiased estimation and better variable selection

