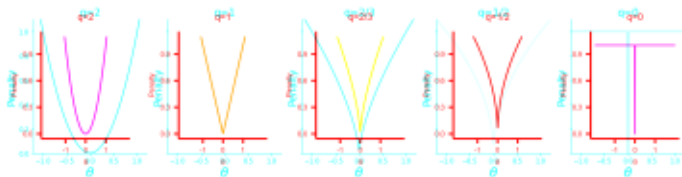


L0 REGULARIZATION

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



- L0 "norm" simply counts the nr of non-zero params
- Induces sparsity more aggressively than L1, but does not shrink
- AIC and BIC are special cases of L0
- L0-regularized risk is not continuous or convex
- NP-hard to optimize; for smaller n and p somewhat tractable, efficient approximations are still current research

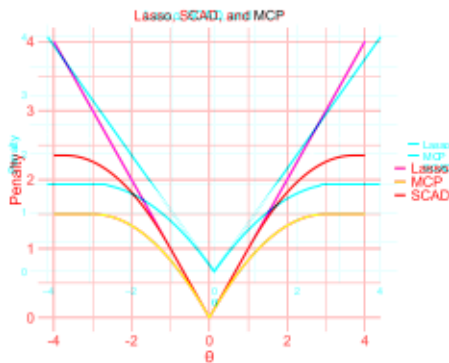
Minimax Concave Penalty:

also non-convex; similar idea as SCAD with $\gamma > 1$

$$MCP(\theta|\lambda, \gamma) = \begin{cases} \lambda|\theta| - \frac{\theta^2}{2\gamma}, & \text{if } |\theta| \leq \gamma\lambda \\ \frac{1}{2}\gamma\lambda^2, & \text{if } |\theta| > \gamma\lambda \end{cases}$$



- As with SCAD, MCP starts by applying same penalization rate as lasso, then smoothly reduces rate to zero as $|\theta| \uparrow$
- Different from SCAD, MCP immediately starts relaxing the penalization rate, while for SCAD rate remains flat until $|\theta| > \lambda$
- Both SCAD and MCP possess oracle property: they can consistently select true model as $n \rightarrow \infty$ while lasso may fail

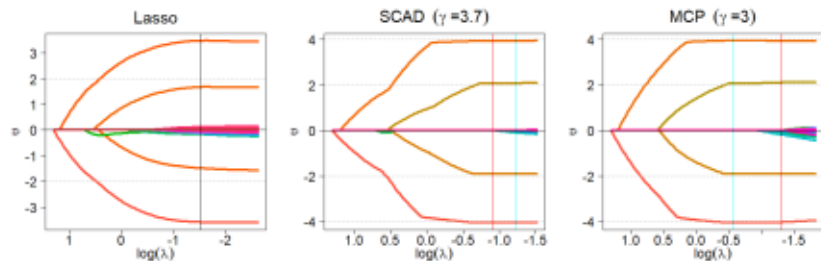


EXAMPLE: COMPARING REGULARIZERS

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

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

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



Vertical lines mark optimal λ from 10CV.

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

