# Introduction to Machine Learning

## Regularization
## Lasso Regression



Effect of L1 Regularization on Linear Model Solutions

**Learning goals**

- Lasso regression / $L$1 penalty
- Know that lasso selects features
- Support recovery

# LASSO REGRESSION

Another shrinkage method is the so-called **lasso regression** (least absolute shrinkage and selection operator), which uses an $L1$ penalty on $\theta$:
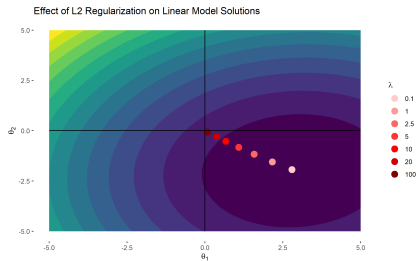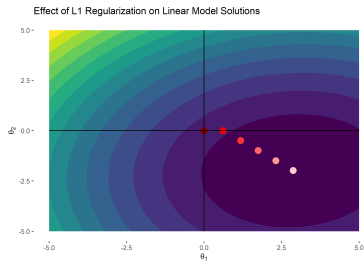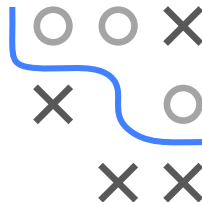
$$\hat{\theta}_{\text{lasso}} = \arg\min_{\theta} \sum_{i=1}^{n} \left( y^{(i)} - \theta^T \mathbf{x}^{(i)} \right)^2 + \lambda \sum_{j=1}^{p} |\theta_j|$$

$$= \arg\min_{\theta} \left( \mathbf{y} - \mathbf{X}\theta \right)^{\top} \left( \mathbf{y} - \mathbf{X}\theta \right) + \lambda \|\theta\|_1$$

Optimization is much harder now. $\mathcal{R}_{\text{reg}}(\theta)$ is still convex, but in general there is no analytical solution and it is non-differentiable.
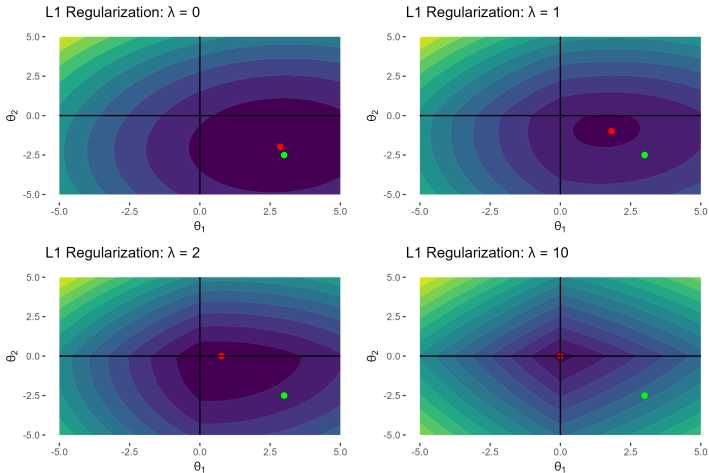
# LASSO REGRESSION / 2

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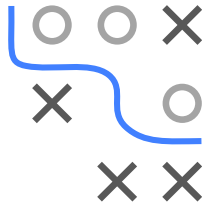$\theta^* = (3, -2)^T$. LHS = $L$1 regularization; RHS = $L$2



With increasing regularization, $\hat{\theta}_{lasso}$ is pulled back to the origin, but
takes a different "route". $\theta_2$ eventually becomes 0!

# LASSO REGRESSION

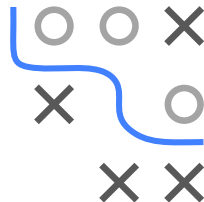Contours of regularized objective for different $\lambda$ values.



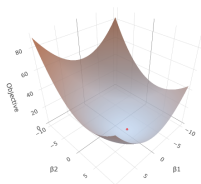Green = true minimizer of the unreg.objective and red = lasso solution.
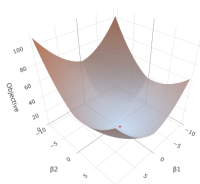
# LASSO REGRESSION / 4

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 $\mathcal{R}_{\text{reg}}$ toward a "basin" (elliptic paraboloid).
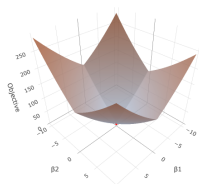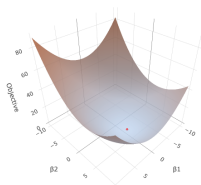
# LASSO REGRESSION / 5

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

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

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.

# *L*1 **AND** *L*2 **REG. WITH ORTHONORMAL DESIGN**

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

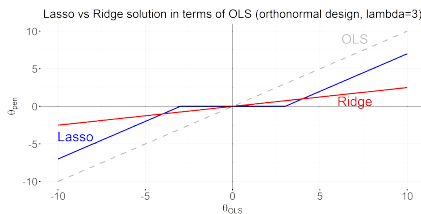$$\hat{\theta}_{\text{lasso}} = \text{sign}(\hat{\theta}_{\text{OLS}})(|\hat{\theta}_{\text{OLS}}| - \lambda)_+ \quad \text{(sparsity)}$$
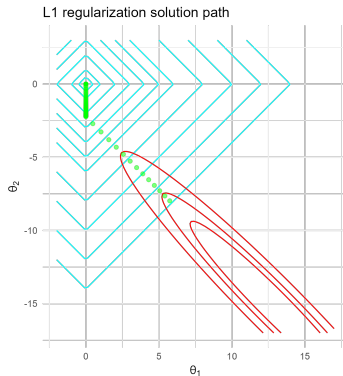
Function $S(\theta, \lambda) := \text{sign}(\theta)(|\theta| - \lambda)_+$ is called **soft thresholding** operator: For $|\theta| \leq \lambda$ it returns 0, whereas params $|\theta| > \lambda$ are shrunken toward 0 by $\lambda$. Comparing this to $\hat{\theta}_{\text{Ridge}}$ under orthonormal design:

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



Lasso vs Ridge solution in terms of OLS (orthonormal design, lambda=3)

# COMPARING SOLUTION PATHS FOR $L1/L2$

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

## SUPPORT RECOVERY OF LASSO  ▸ Zhao and Yu 2006

When can lasso select true support of $\theta$, i.e., only the non-zero parameters?
Can be formalized as sign-consistency:

$$\mathbb{P}\big(\text{sign}(\hat{\theta}) = \text{sign}(\theta)\big) \to 1 \text{ as } n \to \infty \quad (\text{where } \text{sign}(0) := 0)$$

Suppose the true DGP given a partition into subvectors $\theta = (\theta_1, \theta_2)$ is

$$Y = \mathbf{X}\theta + \varepsilon = \mathbf{X}_1\theta_1 + \mathbf{X}_2\theta_2 + \varepsilon \text{ with } \varepsilon \sim (0, \sigma^2 I)$$

and only $\theta_1$ is non-zero. Let $\mathbf{X}_1$ denote the $n \times q$ matrix with the relevant
features and $\mathbf{X}_2$ the matrix of noise features. It can be shown that $\hat{\theta}_{lasso}$ is sign
consistent under an **irrepresentable condition**:

$$|(\mathbf{X}_2^\top \mathbf{X}_1)(\mathbf{X}_1^\top \mathbf{X}_1)^{-1}\text{sign}(\theta_1)| < 1 \text{ (element-wise)}$$

In fact, lasso can only be sign-consistent if this condition holds.
Intuitively, the irrelevant variables in $\mathbf{X}_2$ must not be too correlated with (or
*representable* by) the informative features  ▸ Meinshausen and Yu 2009