

# Regularization

UE de Master 2, AOS1  
Fall 2022

---

S. Rousseau

# General definition

What is regularization? From Goodfellow et al. 2016:

*Regularization is any modification we make to a learning algorithm that is intended to reduce its **generalization error** but not its **training error**.*

- Adding a penalty term to a loss function
- Data augmentation
- Early stopping
- ...

Today's course is focused on the first point

# What is learning

Learning consists in minimizing a training objective

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \hat{L}(f)$$

- $\mathcal{H}$  is the set of admissible classifier/regression function
- $\hat{L}$  is an empirical loss function (computed on a train set)
- $\hat{f}$  is the learnt solution

Most of the time the set  $\mathcal{H}$  is parametrized by a parameter  $\theta \in \Theta$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \hat{L}(\theta)$$

# Application to linear regression

- Set of admissible solutions is linear functions

$$\mathcal{H} = \{\text{linear functions}\}$$

- Linear functions on  $\mathbb{R}^p$  are parametrized by  $\beta \in \mathbb{R}^p$

$$\mathcal{H} = \{\mathbf{x} \mapsto \langle \mathbf{x}, \beta \rangle, \mathbf{x} \in \mathbb{R}^p\}$$

- Training objective is the residual sum of squares (RSS)  
Empirical loss function based on square loss, prediction is  $\langle \mathbf{x}_i, \beta \rangle$ , observed is  $y_i$

$$\hat{L}(\beta) = \text{RSS}(\beta) = \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \beta \rangle)^2$$

## Application to linear regression (2)

- The learning algorithm is then

$$\hat{\beta}^{\text{ols}} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 \quad (\text{ordinary least square})$$

- Compact matrix formulation  $X = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$  and  $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

$$\hat{\beta}^{\text{ols}} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - X\beta\|^2$$

# Regularization

- Unregularized objective

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \hat{L}(\theta) \quad (1)$$

- Regularized objective

$$\hat{\theta}_{\lambda}^{\text{reg}} = \arg \min_{\theta \in \Theta} \hat{L}(\theta) + \lambda \cdot R(\theta)$$

- Training objective
- Tuning parameter
- Regularization term



## Regularization (2)

Regularized objective

$$\hat{\theta}^{\text{reg}} = \arg \min_{\theta \in \Theta} \hat{L}(\theta) + \lambda \cdot R(\theta)$$

- $R(\theta)$  penalizes some  $\theta$ 's
- $\lambda \geq 0$  is the strength of the penalty
  - $\lambda = 0$ : no penalty: regular solution
  - $\lambda \rightarrow +\infty$ , solution tends to  $\arg \min_{\theta \in \Theta} R(\theta)$
  - Some tradeoff has to be found between the two extreme cases

# Ridge regularization

Most simple regularization we can think of:

- We choose  $R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|^2 = \sum_{i=1}^p \theta_i^2$
- Penalizes large parameter: prevents the  $\beta_i$ 's from exploding
- Ridge regularization is then

$$\hat{\boldsymbol{\theta}}_{\lambda}^{\text{ridge}} = \arg \min_{\boldsymbol{\theta} \in \Theta} \hat{L}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2 \quad (\text{ridge regularization})$$

Also known as  $L_2$ -regularization, Tikhonov regularization or weight decay (neural network)



# Application to ridge regression

- Previous linear regression learning algorithm was

$$\hat{\beta}^{\text{ols}} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - X\beta\|^2$$

- Adding the ridge regularizer term yields

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - X\beta\|^2 + \lambda \|\beta\|^2 \quad (\text{ridge regression})$$

$$= \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 + \lambda \sum_{i=1}^p \beta_i^2$$

## Solution to ridge regression

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - X\beta\|^2 + \lambda \|\beta\|^2 \quad (\text{ridge regression})$$

- Define the penalized residual sum of squares (PRSS) as

$$\text{PRSS}(\beta) = \|\mathbf{y} - X\beta\|^2 + \lambda \|\beta\|^2$$

- PRSS is (strictly) convex w.r.t.  $\beta$ : unique solution
- By differentiating w.r.t.  $\beta$  we get

$$\nabla_{\beta} \text{PRSS} = -2X^T(\mathbf{y} - X\beta) + 2\lambda\beta \quad (2)$$

## Solution to ridge regression

- Setting the derivative to zero we finally get

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \left(X^T X + \lambda I_p\right)^{-1} X^T \mathbf{y} \quad (3)$$

- Fitted values are then

$$\hat{\mathbf{y}}^{\text{ridge}} = X \hat{\beta}_{\lambda}^{\text{ridge}} = X \left(X^T X + \lambda I_p\right)^{-1} X^T \mathbf{y}$$

- For  $\lambda = 0$  we have the OLS solution

$$\beta_{\lambda=0}^{\text{ridge}} = \beta^{\text{ols}}$$

$$\hat{\mathbf{y}}^{\text{ols}} = X \left(X^T X\right)^{-1} X^T \mathbf{y}$$

- The intercept (if present) should not be part of the regularizing parameter. Two possible strategies:
  - center the design matrix  $X$  so there is no intercept
  - or we remove the intercept from the regularizing parameter (set  $\beta^*$  as  $\beta$  without  $\beta_0$ )
- The features should be on the **same scale**; unlike linear regression ridge regression predictions are sensitive to features rescaling

# Properties of ridge regression

- Unlike linear regression, there is always a solution (when  $\lambda > 0$ )
  - $X^T X$  is positive **semi-definite**
  - $X^T X + \lambda I_p$  is then positive definite when  $\lambda > 0$  hence is invertible
- It improves the **conditioning** of the problem
- Like linear regression but unlike Lasso, it admits a **closed form solution**

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \left( X^T X + \lambda I_p \right)^{-1} X^T \mathbf{y}$$

- Invariant to rotation: if  $Y = XU^T$  is a rotation of the samples then  $\hat{\beta}_Y = U\hat{\beta}_X$
- Unlike linear regression, both the  $\beta_i$ 's estimate and predictions are **biased**
- The  $\beta_i$ 's estimate are **drawn toward zero** w.r.t the OLS solution
- Might have lower variance than OLS

## $\hat{\beta}_\lambda^{\text{ridge}}$ is biased

Suppose that the design matrix  $X$  is fixed (conditioning on it)

Linear case:

$$\begin{aligned}\mathbb{E}(\hat{\beta}^{\text{ols}}) &= \mathbb{E}\left(\left(X^T X\right)^{-1} X^T \mathbf{y}\right) \\ &= \left(X^T X\right)^{-1} X^T \mathbb{E}(\mathbf{y}) \\ &= \left(X^T X\right)^{-1} X^T X \beta \\ &= \beta\end{aligned}$$

Unbiased!

Ridge case:

$$\begin{aligned}\mathbb{E}(\hat{\beta}_\lambda^{\text{ridge}}) &= \mathbb{E}\left(\left(X^T X + \lambda I_p\right)^{-1} X^T \mathbf{y}\right) \\ &= \left(X^T X + \lambda I_p\right)^{-1} X^T X \beta \\ &= \left(\left(X^T X\right)^{-1} \left(X^T X + \lambda I_p\right)\right)^{-1} \beta \\ &= \left(I_p + \lambda \left(X^T X\right)^{-1}\right)^{-1} \beta \neq \beta\end{aligned}$$

Biased!

## Data augmentation interpretation

- Let's rewrite the PRSS

$$\begin{aligned}\text{PRSS}(\boldsymbol{\beta}) &= \|\mathbf{y} - X\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \\ &= \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle)^2 + \lambda \sum_{i=1}^p \beta_i^2 \\ &= \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle)^2 + \sum_{i=1}^p \left( 0 - \langle \sqrt{\lambda} \mathbf{e}_i, \boldsymbol{\beta} \rangle \right)^2\end{aligned}$$

- Same as adding  $p$  extra samples in addition to the  $n$   $\mathbf{x}_i$ 's
- Additional samples and observations are  $(\sqrt{\lambda} \mathbf{e}_i, 0)$  for  $i = 1, \dots, p$
- Same as adding an observation on each axis that is zero

## Data augmentation interpretation (cont'd)

- Switching back to matrix form we define

$$X_{\lambda} = \begin{pmatrix} X & & & \\ \sqrt{\lambda} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda} & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \sqrt{\lambda} \end{pmatrix} \quad \text{and} \quad \mathbf{y}' = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

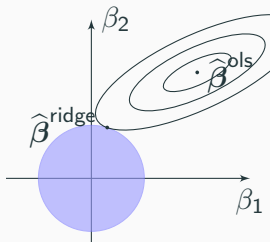
- And the PRSS can now be written

$$\|\mathbf{y}' - X_{\lambda}\boldsymbol{\beta}\|^2 \quad \text{s.t.} \quad X_{\lambda} = \begin{pmatrix} X \\ \sqrt{\lambda}I_p \end{pmatrix} \quad \mathbf{y}' = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}$$



## Geometric interpretation $\hat{\beta}^{\text{ridge}}$

- Suppose 2-dimensional case ( $p = 2$ ),  $\hat{\beta}^{\text{ols}}$  is the OLS solution
- Ellipses are level line of the RSS:  $\|\mathbf{y} - X\beta\|^2$
- A solution for some  $\lambda$  is at the intersection of the  $L_2$  ball and a level line
- Whatever the form of ellipses, ridge solution is systematically drawn toward zero

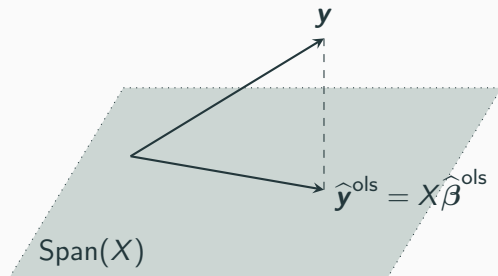


# Geometric interpretation of $\hat{\mathbf{y}}^{\text{ols}}$

First see the OLS case

- Let  $X = USV^T$  the SVD of  $X$
- Matrix  $U$  gathers the (unit) principal components  $\mathbf{u}_1, \dots, \mathbf{u}_k$
- Ordinary least squares orthogonally projects  $\mathbf{y}$  onto the space spanned by the columns of  $X$ :

$$\begin{aligned}\hat{\mathbf{y}}^{\text{ols}} &= X(X^T X)^{-1} X^T \mathbf{y} = U U^T \mathbf{y} \\ &= \sum_{i=1}^p (\mathbf{u}_i^T \mathbf{y}) \mathbf{u}_i\end{aligned}$$



## Geometric interpretation of $\hat{\mathbf{y}}^{\text{ridge}}$ (cont'd)

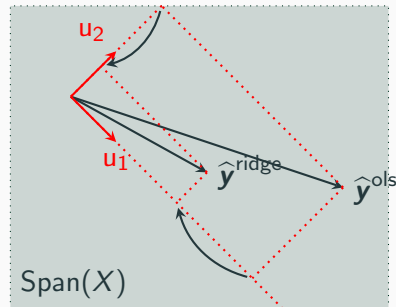
- Ridge regression is doing the same thing plus an additional “shrinking” step

$$\begin{aligned}\hat{\mathbf{y}}^{\text{ridge}} &= \mathbf{X} \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{U} \left( \mathbf{S} (\mathbf{S}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{S} \right) \mathbf{U}^T \mathbf{y} \\ &= \sum_{i=1}^p \left( \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \mathbf{u}_i^T \mathbf{y} \right) \mathbf{u}_i\end{aligned}$$

- Coordinates are now shrunk towards zero since  $\frac{\sigma_i^2}{\sigma_i^2 + \lambda} < 1$
- Remember that the  $\mathbf{u}_i$  are here the (unit) principal components of  $\mathbf{X}$
- The lesser the variance of the principal component is the greater it is shrunk
- Smooth version of PCA followed by linear regression

## Geometric interpretation of $\hat{\mathbf{y}}^{\text{ridge}}$ (cont'd)

- Coordinate of  $\hat{\mathbf{y}}^{\text{ols}}$  along  $\mathbf{u}_1$  is shrunk by  $\frac{\sigma_1^2}{\sigma_1^2 + \lambda}$
- Coordinate of  $\hat{\mathbf{y}}^{\text{ols}}$  along  $\mathbf{u}_2$  is shrunk by  $\frac{\sigma_2^2}{\sigma_2^2 + \lambda}$



## Geometric interpretation of $\hat{\beta}_\lambda^{\text{ridge}}$

What is the link between the unknown parameter  $\beta$  and its ridge estimate  $\hat{\beta}_\lambda^{\text{ridge}}$ ?

- Let  $X = USV^T$  the SVD of  $X$
- Recall that  $V$  gathers the principal directions (new basis of representation)
- Let's look at  $\hat{\beta}_\lambda^{\text{ridge}}$  in basis  $V$ , we can show that

$$V^T \hat{\beta}_\lambda^{\text{ridge}} = (I_p + \lambda S^{-2})^{-1} V^T \hat{\beta}^{\text{ols}}$$

In the basis defined by  $V$ ,  $\beta$  is shrunk

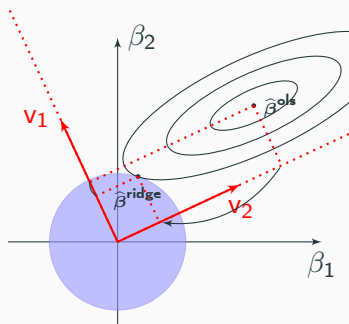
- If we look at the  $i$ th coordinate

$$\left( V^T \hat{\beta}_\lambda^{\text{ridge}} \right)_i = \frac{\sigma_i^2}{\lambda + \sigma_i^2} \left( V^T \hat{\beta}^{\text{ols}} \right)_i$$

# Geometric interpretation of $\hat{\beta}_\lambda^{\text{ridge}}$

$$\left(V^T \hat{\beta}_\lambda^{\text{ridge}}\right)_i = \frac{\sigma_i^2}{\lambda + \sigma_i^2} \left(V^T \hat{\beta}^{\text{ols}}\right)_i$$

- Coordinate of  $\hat{\beta}_\lambda^{\text{ridge}}$  along  $v_1$  is shrunk by  $\frac{\sigma_1^2}{\sigma_1^2 + \lambda}$
- Coordinate of  $\hat{\beta}_\lambda^{\text{ridge}}$  along  $v_2$  is shrunk by  $\frac{\sigma_2^2}{\sigma_2^2 + \lambda}$



## Gradient descent interpretation

- General gradient descent update of step  $\eta$

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \eta \nabla L(\boldsymbol{\theta}^k),$$

- Gradient descent update for OLS ( $L = \text{RSS}$ )

$$\boldsymbol{\beta}^{k+1} = \boldsymbol{\beta}^k - \eta \nabla \text{RSS}(\boldsymbol{\beta}^k) \quad (4)$$

- Gradient descent update for ridge regression

$$\boldsymbol{\beta}^{k+1} = \boldsymbol{\beta}^k - \eta \nabla \text{PRSS}(\boldsymbol{\beta}^k)$$

- Which writes

$$\boldsymbol{\beta}^{k+1} = (1 - 2\eta\lambda)\boldsymbol{\beta}^k - \eta \nabla \text{RSS}(\boldsymbol{\beta}^k) \quad (5)$$

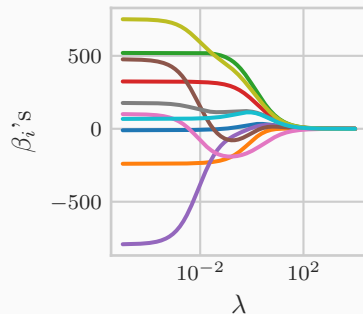
The update use a shrunk version of  $\boldsymbol{\beta}^k$

## Effect on the $\beta_i$ 's: the regularization path

**Regularization path** : plot of  $\beta_i$ 's against regularization parameter  $\lambda$ .

Here for some centered dataset with 10 covariates:

- Linear regression estimates at  $\lambda = 0$
- When regularization is too strong, we fit the constant function
- $\beta_i$ 's are shrunk smoothly towards 0 as  $\lambda$  increases
- All the  $\beta_i$ 's might not be non-increasing but  $\|\beta\|$  is
- The  $\beta_i$ 's are never zero

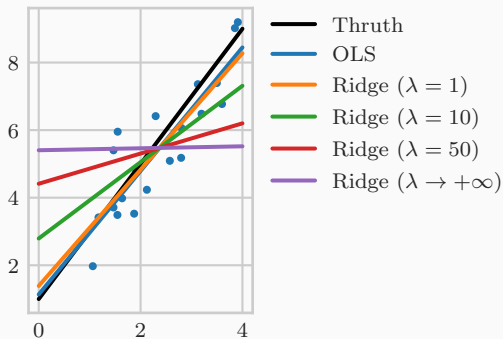




## Effect on the fitted line

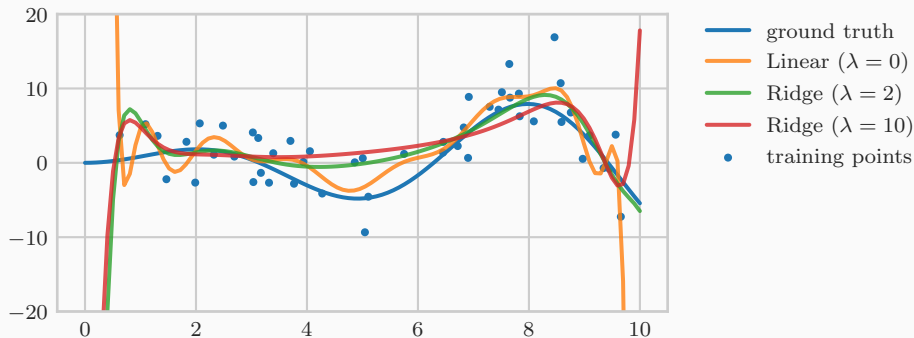
In the case of simple linear regression ( $p = 1$ )

- Regularizing is like adding the point  $(\sqrt{\lambda}, 0)$
- As  $\lambda$  goes to infinity, we fit a constant function



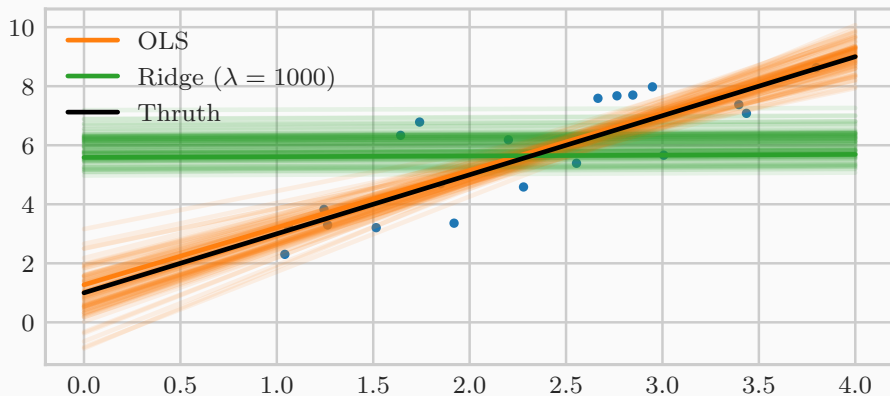
## Effect on fitted curve

Adding polynomial features (degree 15):  $X_i, X_i^2, \dots, X_i^{15}$



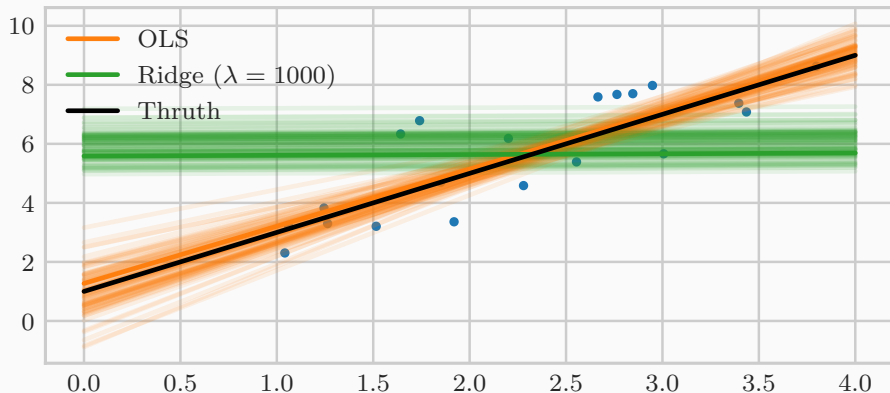
## Effect on the $\beta_i$ 's estimates: bias–variance tradeoff

The slope  $\beta_1$  is biased but variance of estimations is smaller



## Effect on predictions: bias–variance tradeoff

Predictions at 0 (for example) are biased but less spread out



# Lasso regularization

## Other famous regularization

- We choose  $R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_{i=1}^p |\theta_i|$
- Penalizes large parameter: prevents the  $\beta_i$ 's from exploding
- Lasso regularization is then

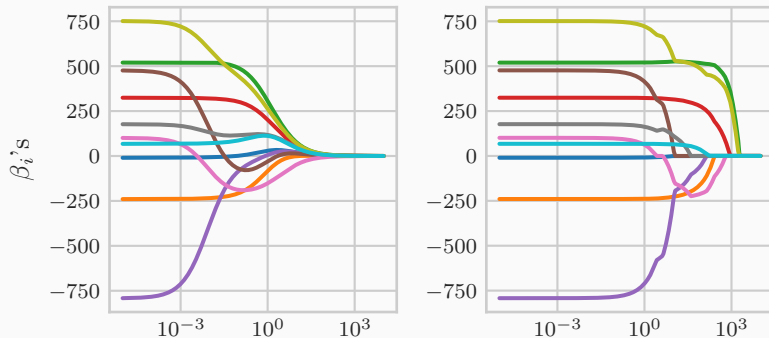
$$\hat{\boldsymbol{\theta}}_{\lambda}^{\text{lasso}} = \arg \min_{\boldsymbol{\theta} \in \Theta} \hat{L}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_1 \quad (\text{Lasso regularization})$$

- Lasso linear regression

$$\hat{\boldsymbol{\beta}}_{\lambda}^{\text{lasso}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{y} - X\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 \quad (\text{Lasso regression})$$

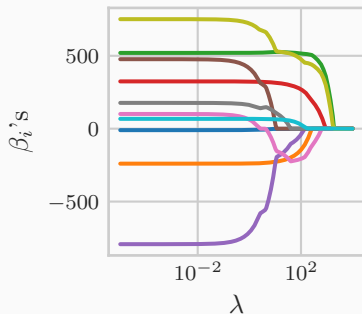
# Sparsity promoting property

Compare the coefficients  $\beta_i$ 's with ridge and lasso regularization



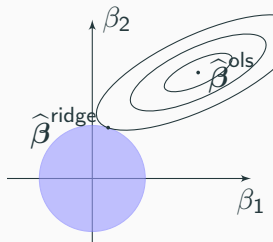
## Effect on the $\beta_i$ 's: the regularization path

- Linear regression estimates at  $\lambda = 0$
- Piecewise linear regularization path
- $\beta_i$ 's are shrunk as  $\lambda$  increases
- All the  $\beta_i$ 's are shrunk to exactly zero at some point: **sparsity promoting effect**
- When regularization is too strong, we fit the constant function

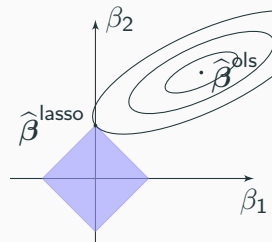


## Explaining the sparsity property

- $\hat{\beta}^{\text{ols}}$  is the ordinary least square solution



Ridge can be anywhere on the  $L_2$  ball



Lasso solution lies on edge of  $L_1$  ball  
(sparse solution)



# Geometric interpretation of $\hat{\beta}_\lambda^{\text{lasso}}$

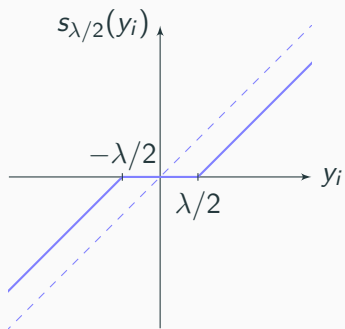
For simplicity, suppose that  $n = p$  and  $X$  is the identity matrix

- OLS solution is:  $\hat{\beta}^{\text{ols}} = \mathbf{y}$
- Lasso regularization reads:  $\hat{\beta}_\lambda^{\text{lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \beta\|^2 + \lambda \|\beta\|_1$

$$\left(\hat{\beta}_\lambda^{\text{lasso}}\right)_i = \arg \min_{\beta_i \in \mathbb{R}} (y_i - \beta_i)^2 + \lambda |\beta_i|$$

$$\left(\hat{\beta}_\lambda^{\text{lasso}}\right)_i = \begin{cases} \max(y_i - \lambda/2, 0) & \text{if } y_i \geq 0 \\ \max(y_i + \lambda/2, 0) & \text{if } y_i < 0 \end{cases}$$

(soft thresholding)



## Gradient descent interpretation

- Using subgradient to differentiate  $\|\cdot\|_1$

$$\nabla_{\beta} \|\beta\|_1 = \text{sign}(\beta)$$

- Gradient descent update for lasso regression

$$\beta^{k+1} = \beta^k - \eta \nabla \text{RSS}(\beta^k) - \eta \lambda \text{sign}(\beta^k)$$

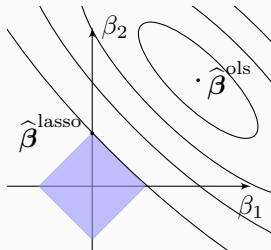
- Shrink the  $\beta_i$ 's regardless of the  $\beta_i$ 's magnitude

# Lasso properties

- No closed form solution
- Convex problem
- Feature selection ability
- Biased predictions and parameter estimate
- Might be unstable if highly correlated variables

## Why elastic-net?

- Ridge regression is not selecting variables (all the  $\beta_i$ 's are nonzero)
- Lasso regression is but in an unstable way
- Small changes in  $X$  might lead to entirely different set of selected predictors

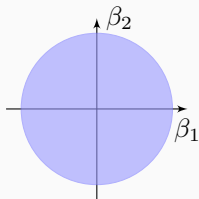


# Elastic net

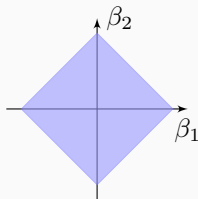
Mixing the two strategies

$$\hat{\boldsymbol{\theta}}_{\lambda, \alpha}^{\text{elastic}} = \arg \min_{\boldsymbol{\theta} \in \Theta} \hat{L}(\boldsymbol{\theta}) + \lambda \left( \alpha \|\boldsymbol{\theta}\|_1 + (1 - \alpha) \|\boldsymbol{\theta}\|_2^2 \right) \quad (\text{Elastic net regularization})$$

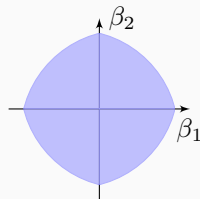
- $\lambda$  is the regularizing parameter
- $\alpha$  controls the balance between  $L_1$  and  $L_2$  regularizing terms



(a)  $L_2$ : no selection



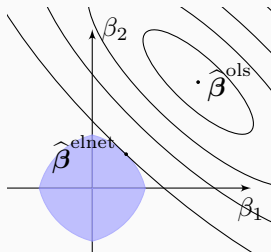
(b)  $L_1$ : unstable selection



(c) Elastic net: stable selection

## Explaining the stability of elastic net regularization

- Corners are still sharp: elastic net is still encouraging sparsity
- Elastic net ball is also round (4 portions of (big) circles): stable when some variables are strongly correlated



# Ridge/Lasso/Elastic net regression in Python and Scikit-Learn

- Import the PCA module

```
from sklearn.linear_model import LinearRegression, Ridge, Lasso
```

- Instantiate (no parameter), fit and predict

```
lr = LinearRegression()  
lr.fit(X, y)  
res = lr.predict(new_X)
```

- Instantiate with tuning parameter alpha

```
lr = Ridge(alpha=1.0)  
lr.fit(X, y)  
res = lr.predict(new_X)
```

- [1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning*. Vol. 1. Springer series in statistics New York, 2001.
- [2] Ian Goodfellow et al. *Deep Learning*. Vol. 1. MIT press Cambridge, 2016.