

Regularization

Constraint it

Constraint Model \rightsquigarrow Solve overfitting

- For polynomial \Rightarrow constraint degree
- For linear \Rightarrow actually constraint Weights

Table of Contents :

- ① Ridge Regression
 - ② Lasso Regression
 - ③ Elastic Net
- Three different way to constraint weights

① Ridge Regression

(Tikhonov Regularization)

- Regression + Regularizations

with L_2 norm

- let $\vec{w} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_N \end{bmatrix}$ $\|\vec{w}\|_2^2 = \theta_1^2 + \theta_2^2 + \dots$

- Cost function after Regularization

$$J(\vec{\theta}) = \text{MSE}(\theta) + \frac{\alpha}{2} \sum_{i=1}^N \theta_i^2$$

\downarrow
 $\|\vec{w}\|_2^2$

- **Notes** :

- 1 - it penalty large weights of features (Fit train set with small weights as possible $\alpha \rightarrow 0$)
 - \rightarrow large α will lead model to be \Rightarrow flat line around average
 - \rightarrow small α will lead model \simeq Normal linear Regress.

2 - Bias term not added to regularization

3 - use only Regularized performance measure in training

NOTE

It is quite common for the cost function used during training to be different from the performance measure used for testing. Apart from regularization, another reason they might be different is that a good training cost function should have optimization-friendly derivatives, while the performance measure used for testing should be as close as possible to the final objective. For example, classifiers are often trained using a cost function such as the log loss (discussed in a moment) but evaluated using precision/recall.

For training :

1 - Closed Form

before : $\hat{\theta} = (X^T X)^{-1} X^T y$

after : $\hat{\theta} = (X^T X + \alpha \mathbf{A})^{-1} X^T y$

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}_{(n+1) \times (n+1)}$$

2 - Gradient descent

$$J(\theta) = \text{MSE}(\theta) + \frac{\alpha}{2} \sum_{i=1}^n \theta_i^2$$

$$\nabla J = \nabla \text{MSE}(\theta) + \alpha \vec{w}$$

Just derive $\frac{d}{d\theta_2} (\frac{\alpha}{2} \theta_2^2) = \alpha \theta_2$

Observation

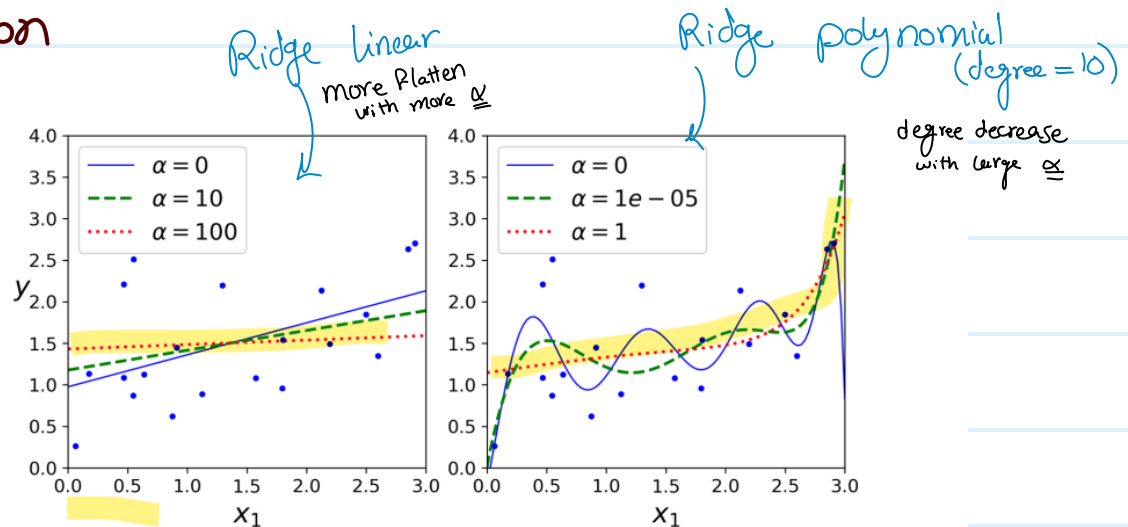


Figure 4-17. A linear model (left) and a polynomial model (right), both with various levels of Ridge regularization

with " α large" \leadsto underfit
" α small" \leadsto Overfit

Must Rescale // implementation in Note Book

② Lasso Regression

- ① — least absolute shrinkage & selection operator
- ② — add L_1 norm instead of $\frac{1}{2} L_2$ norm
- ③ — need less " α " values than ridge

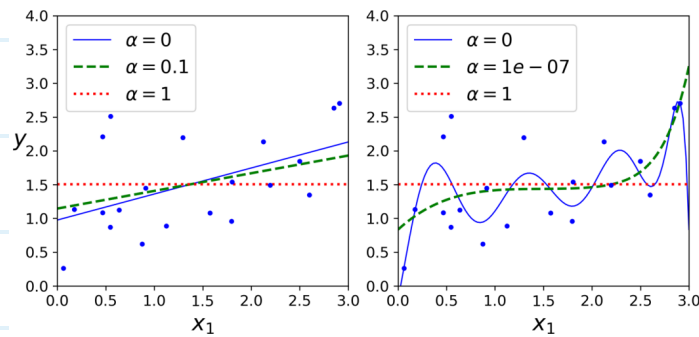


Figure 4-18. A linear model (left) and a polynomial model (right), both using various levels of Lasso regularization

$$J(\theta) = \text{MSE} + \alpha \sum_{i=1}^n \theta_i$$

\vec{g} instead of ∇J as it's not differentiable when " $\theta_i = 0$ "

$$\vec{g} = \nabla_{\theta} \text{MSE}(\theta) + \alpha \begin{pmatrix} \text{sign}(\theta_1) \\ \text{sign}(\theta_2) \\ \vdots \\ \text{sign}(\theta_n) \end{pmatrix}$$

$$\text{sign}(x) = \begin{cases} +1 \\ 0 \\ -1 \end{cases}$$

\vec{g} = subgradient vector

⑤ tend to eliminate \Rightarrow the weights of least important features

↳ output sparse model "few non zero features"

↳ automatically feature selection

how? or why?

⑩

\Rightarrow implementation in Note base

⑥ Why & How?

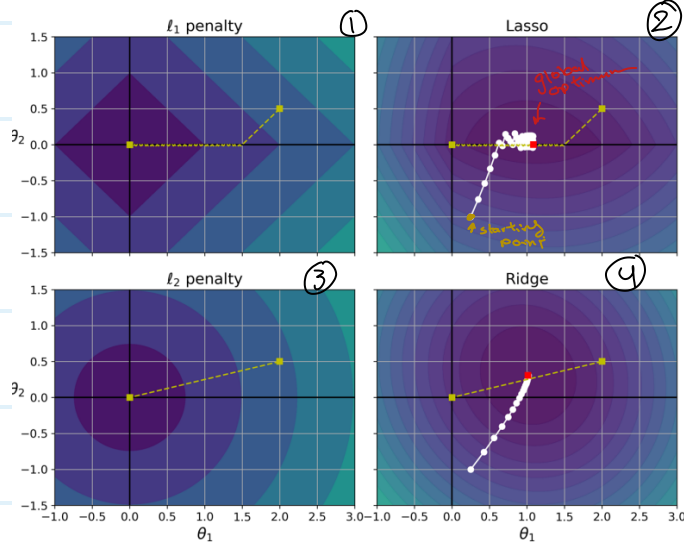


Figure 4-19. Lasso versus Ridge regularization

① L₁ penalty

$$\text{Loss} = J(\vec{\theta}) = \sum_{i=1}^n \theta_i^2$$

$$\nabla_{\theta} J(\theta) = \begin{pmatrix} \text{sign}(\theta_1) \\ \vdots \end{pmatrix}$$

$$\theta_i = \theta_i - \eta \text{sign}(\theta_i) \quad \leftarrow \text{gradient descent}$$

→ it decrease all features linearly till reach all zero (or close as it update $(\pm \eta, \text{zero})$)

→ it start @ $\vec{\theta} = \begin{pmatrix} 2.0 \\ 0.5 \end{pmatrix}$
 same as optimal parameters for unregularized gradient "MSE"

② Lasso

$$\vec{J}(\vec{\theta}) = \text{MSE}(\vec{\theta}) + \alpha \sum_{i=1}^n \theta_i$$

$$\nabla J(\vec{\theta}) = \nabla \text{MSE}(\vec{\theta}) + \alpha \begin{pmatrix} \text{sign}(\theta_1) \\ \vdots \\ \text{sign}(\theta_n) \end{pmatrix}$$

$$\theta_i = \theta_i - \eta (\nabla_{\theta} \text{MSE}(\vec{\theta}) + \alpha (!))$$

$$= \theta_i - \eta (\text{unregularized (MSE)} + \alpha L_1)$$

if we decrease " α " red point move Right to unregularized global optimum

if we increase " α " red point move Left to regularized " "

Note : it keep Bouncing around global optimum in Both ①, ② images

→ as it update By $\pm \eta, \alpha$ or zero

→ One Solution's to ^{gradually} decrease η

→ it will bounce around But Steps will get smaller till Converge

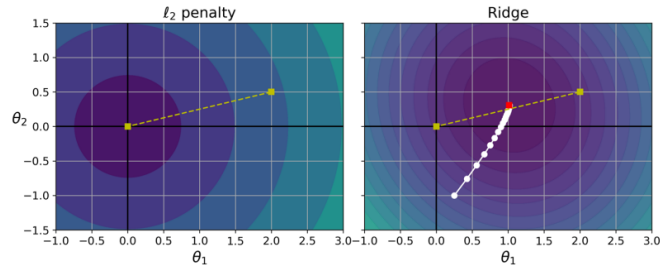
③ l_2 penalty , ④ Ridge

$$J_{l_2} = \frac{\alpha}{2} \theta_i^2$$

$$= \frac{\alpha}{2} \cdot \text{Squared distance of parameters from origin}$$

= like circle

→ gradient descent just take straight path to origin



Main two Difference → gradients get smaller as it approach global minimum
 → making it Converge ⇒ no Bouncing around

→ the optimal parameters "red point" come close to zero as we increase " α " But Never get eliminated (true zero)

③ Elastic Net

- middle ground Between $\xrightarrow{\text{Ridge}}$
 $\xrightarrow{\text{Lasso}}$

$$J(\theta) = \text{MSE}(\theta) + r \alpha \sum_{i=1}^n |\theta_i| + (1-r) \frac{\alpha}{2} \sum_{i=1}^n \theta_i^2$$

$$r = \begin{cases} 0 \rightarrow \text{Ridge} \\ 1 \rightarrow \text{Lasso} \end{cases}$$

\longrightarrow implementation in Notebook

* Frame Work to Choose From

① Ridge is Good Default

② Elastic Net

\rightarrow When Suspecting only Few Features are useful
 \rightarrow preferred over "Lasso" \rightarrow as Lasso perform Randomly when

① $\# \text{Features} > \# \text{samples}$

② Several Features are strong correlated

③ Lasso

④ Plain Linear Regression

\rightarrow it always prefer to have some \uparrow Regularization
 or little bit of

④ Early Stopping

- Stop training when Validation Reach its minimum
 & Before starting to go up again "Before overfit"

- called "Beautiful Free lunch"

- For "SGD", "Mini-Batch GD" the curve will be
Bumpier & will be hard to determine "min point"
 ↳ Solution is to wait some time till be sure
 then Roll-Back

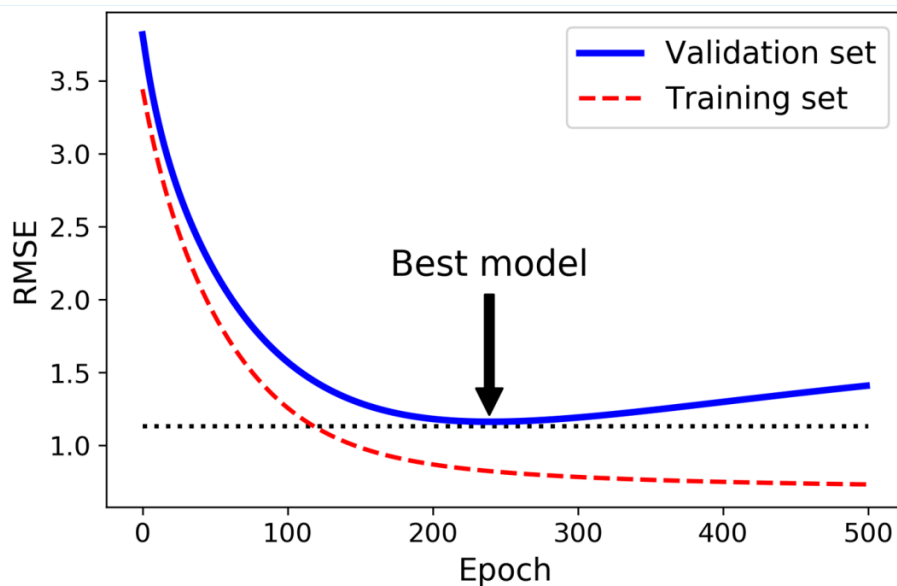


Figure 4-20. Early stopping regularization

→ implementation
 in Notebook