

## Ridge Regression

LS loss fn ① + ② + L2 penalized mean loss ③

$$\text{find } w^* = \underset{w}{\operatorname{argmin}} \|Xw - \vec{y}\|_2^2 + \lambda \|w\|_2^2 = J(w)$$

- $X$  has fictitious dimension
- don't penalize bias term

regularization / penalty term

don't count bias term

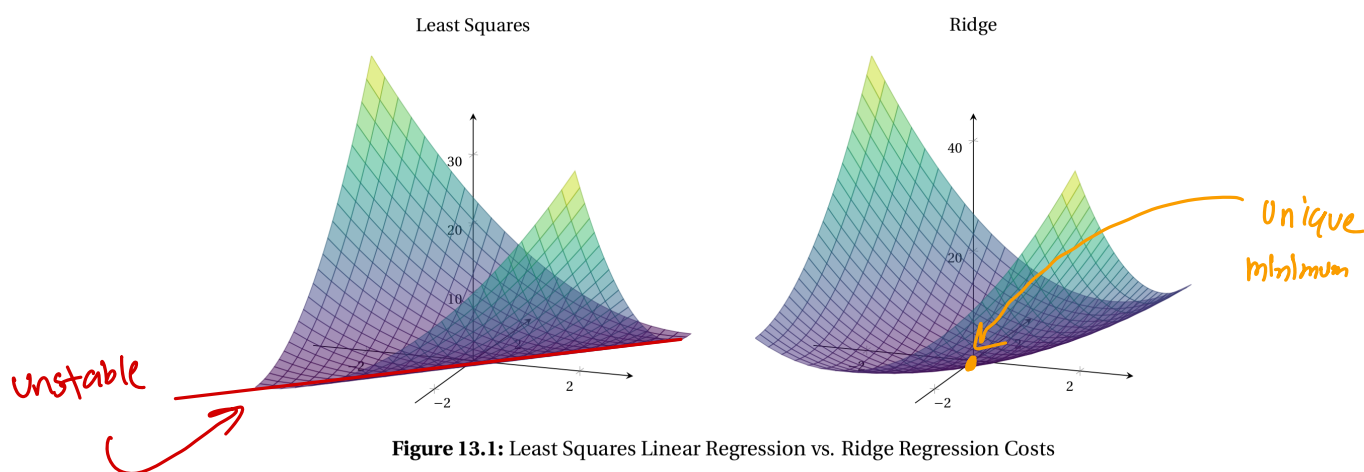
Why regularization?

Shrinkage: make weights small

Why shrinkage?

- guarantees PD normal equations  
↳ always unique solution

e.g. when  $d > n$ ,



"ill-posed" problem: not just 1 solution

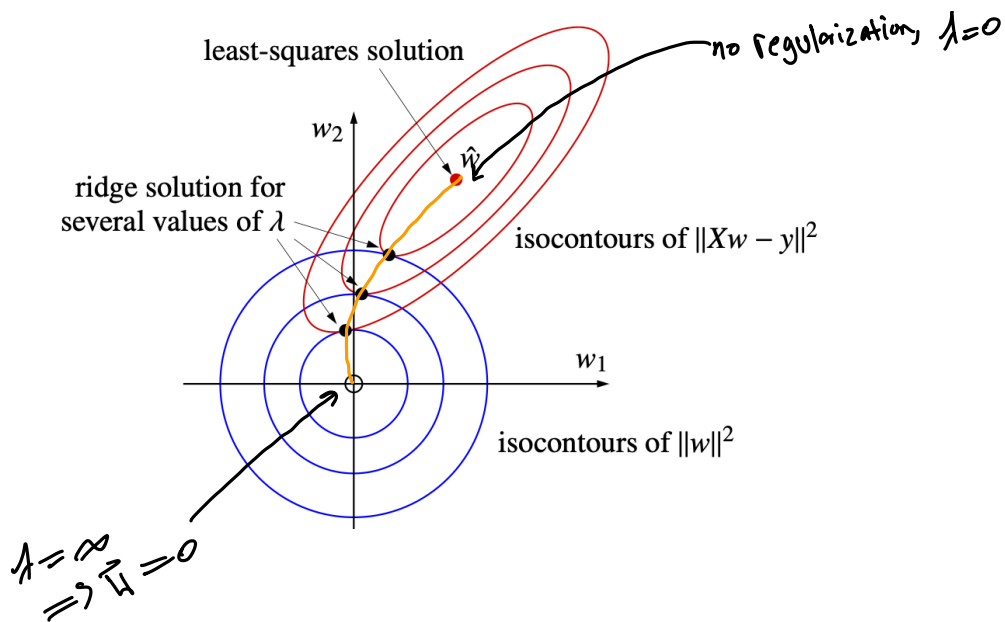
- regularization reduces overfitting  
↳ decreases Variance

Why does this reduce variance?

Imagine:  $500x_1 - 500x_2$  is line of best fit for well-separated pts,  $y_i \in [0, 1]$

↳ large coeffs → instability  
↳ small changes in  $X$  → big change in  $y$

↳ Solution: penalize large weights



Validation to find best  $\lambda$

Set  $\nabla J(w)$  to 0:

$$(X^T X + \lambda I)^{-1} X^T y$$

last diag. entry is 0

$$\Rightarrow \hat{w} = (X^T X + \lambda I)^{-1} X^T y$$

$$\Rightarrow h(z) = \hat{w}^T z$$

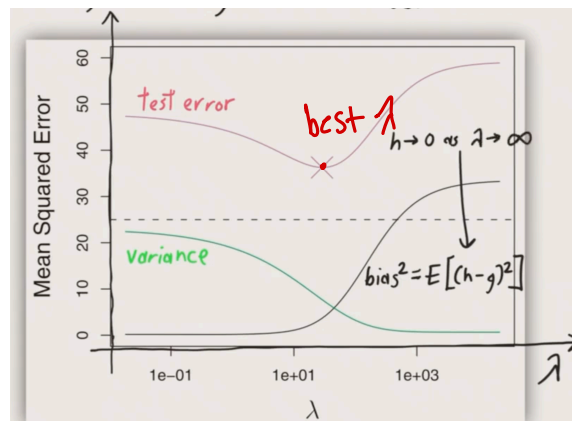
Increasing  $\lambda \rightarrow$  more regularization  $\rightarrow$  smaller  $\|\hat{w}\|_2^2$  & vice versa

recall data model:  $\vec{y} = X\vec{v} + \vec{\epsilon}$ ,  $\vec{\epsilon}$  is noise

$$\text{Var}(\text{ridge regr}) = \text{Var}(\vec{\epsilon}^T (X^T X + \lambda I)^{-1} X^T \vec{e})$$

$\lambda \uparrow$ ,  $\text{Var} \downarrow$

As  $\lambda \rightarrow \infty$ ,  $\text{Var}(\cdot) \rightarrow 0$ , bias  $\uparrow$



- tune  $\lambda$  w/ validation or CV

- Ideally, features should be normalized to have same variance

alternatives: use asymmetric penalty by replacing  $I'$  w/ another diagonal matrix

\* In polynomial regression, different weights need to be penalized differently

### Bayesian Justification for Ridge Regression

- Assign a prior probability on  $W$ :  $W \sim N(0, \sigma^2)$

→ true weights come from this distr.

- apply MLE to maximize posterior probability

Bayes's thm: posterior  $f(W | X, Y) = \frac{f(Y | X, W) \cdot \text{prior } f(W)}{f(Y | X)} = \frac{\mathcal{L}(W) f(W)}{f(Y | X)}$

$W, Y$  RVs,  $X$  is a constant

→ maximize log posterior:  $\ln \mathcal{L}(W) + \ln f(W) - \text{constant}$

$$= -\text{const} \|XW - \vec{y}\|_2^2 - \text{const} \|W'\|^2$$

→ minimize  $\|XW - \vec{y}\|_2^2 + \lambda \|W'\|^2$

### Feature Subset Selection

Recall: all features increase variance, not all features reduce bias

Idea: identify which features we can get rid of

↳ Set their weights to 0

↳ less overfitting, smaller test error

↳ Inference. Simpler models convey interpretable wisdom

- Useful in all classification & regression methods.

- Sometimes it's hard: different features can be redundant

Algo: Best Subset Selection. 1) Try all  $2^d - 1$  nonempty subsets of features  
2) choose best classifier w/ validation  
- slow

Heuristics: 1) forward stepwise selection

- start w/ null model ( $0$  features)

doesn't always  
pick best features

- repeatedly add best feature until validation errors start increasing (due to overfitting) instead of decreasing
- at each outer iteration, inner loop tries every feature & chooses best validation. Requires  $d^2$  models trained (better than  $O(2^d)$ )

e.g. won't find best model w/ 2 features if neither feature yields best 1D model

Heuristic 2) backward step selection

best choice if  
most features  
are helpful

- start w/ all  $d$  features
- repeatedly remove features whose removal reduces validation error
- also  $O(d^2)$

Lasso (Robert Tibshirani)

"least absolute shrinkage & selection operator"

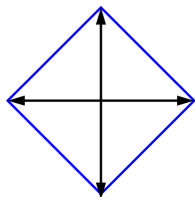
- regression w/  $L_1$  penalty

find  $\vec{w}$  that minimizes:  $J(w) = \|X\vec{w} - \vec{y}\|^2 - \lambda \|\vec{w}\|_1$

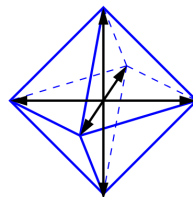
Recall ridge regression: isosurfaces of  $\|\vec{w}\|^2$  are hyperspheres

- those of  $L_1$  are cross-polytopes

- Unit cross-polytope is convex hull of unit coordinate vectors:



$$\|w\|_1 = 1$$



3d  $\rightarrow$  2d slice

