

Lasso Regression Regularization for feature Selection

Why might you want to perform feature selection?

Efficiency:

- If $\text{size}(w) = 100B$, each prediction is expensive
- If w is sparse, computation only depends on # of non-zeros

many zeros

$$\boxed{\quad} = \boxed{\quad \quad \quad \quad \quad}$$

$$\boxed{\quad \quad \quad \quad \quad}$$

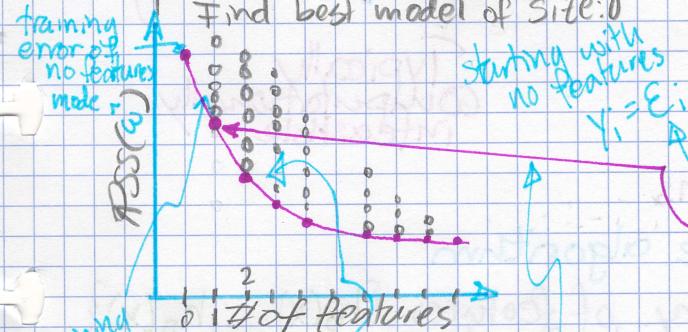
$$y_i = \sum_{\hat{w}_j \neq 0} \hat{w}_j h_j(x_i)$$

Interpretability:

- Which features are relevant for predictions?

Option 1: All Subsets

Find best model of size: 0



training error of model with just one feature

bed as feature

Search over all combinations of 2 features

bed
bath

best fitting model with only one feature

Not necessarily nested

- # bedrooms
- # bathrooms
- sq. ft living
- sq. ft lot
- Floors
- year built
- year renovated
- waterfront

Best model of size K need not contain features of best model of size K-1

Choosing model complexity?

Option 1: Assess on validation set

Option 2: Cross validation

Option 3+: Other metrics for penalizing model complexity like BIC

Complexity of "all subsets"

How many models were evaluated?
- each indexed by features included

$$y_i = \epsilon_i$$

[000..000]

$$y_i = w_0 h_0(x_i) + \epsilon_i$$

[100..000]

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \epsilon_i$$

[010..000]

:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i$$

:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i$$

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

[010..010]

[001..001]

[000..000]

[111..111]

[101..101]

[011..011]

[001..001]

[000..000]

[110..110]

[100..100]

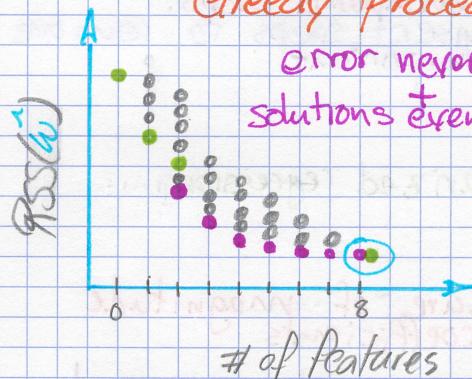
[010..010]

[001..001]

[000..000]

Greedy Procedure

error never increases
solutions eventually meet



- # bedrooms
- # bathrooms
- **Sq ft living** forced to keep in model
- Sq ft lot
- floors
- year built
- year renovated
- waterfront

When do we stop?

When training error is low enough?
NO!

When test error is low enough?
NO!

Use validation set or cross validation!

Complexity of forward stepwise

How many models were evaluated

- 1st step D models
- 2nd step D-1 models
- 3rd step D-2 models
- ...

How many steps?

- Depends

- At most D steps (to full model)

$$\left\{ \begin{array}{l} O(D^2) \ll 2^D \\ \text{for large } D \end{array} \right.$$

Other greedy algorithms

Instead of starting from simple model and always growing.

Backward stepwise: Start with full model and iteratively remove features least useful to fit

Combining forward and backward steps:

In forward algorithm insert steps to remove features no longer as important

Option 3: Regularize

Ridge regression: L₂ regularized regression

Total cost =

measure of fit + measure of magnitude of coefficients

$$\text{RSS}(\mathbf{w}) + \|\mathbf{w}\|_2^2 = w_0^2 + \dots + w_d^2$$

Encourages small weights but not exactly 0

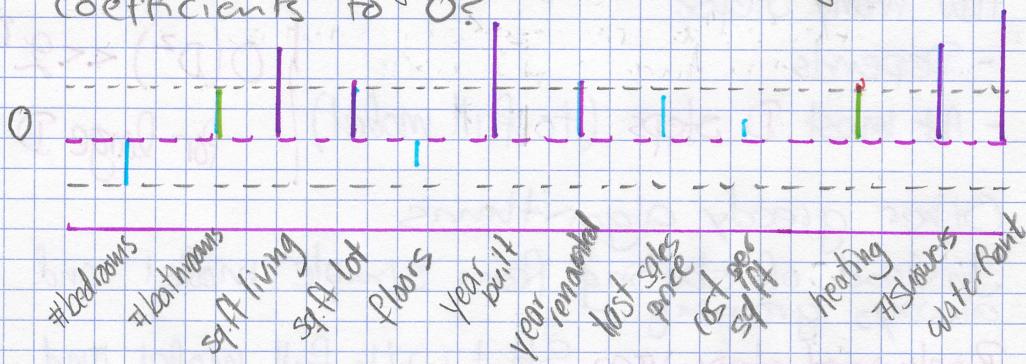
Using regularization for feature selection

Instead of searching over a discrete set of solutions, can we use regularization?

- Start with full model (all possible features)
- "Shrink" some coefficients exactly to 0
• i.e. knock out certain features
- Non-zero coefficients indicate "selected" features

Thresholding ridge coefficients?

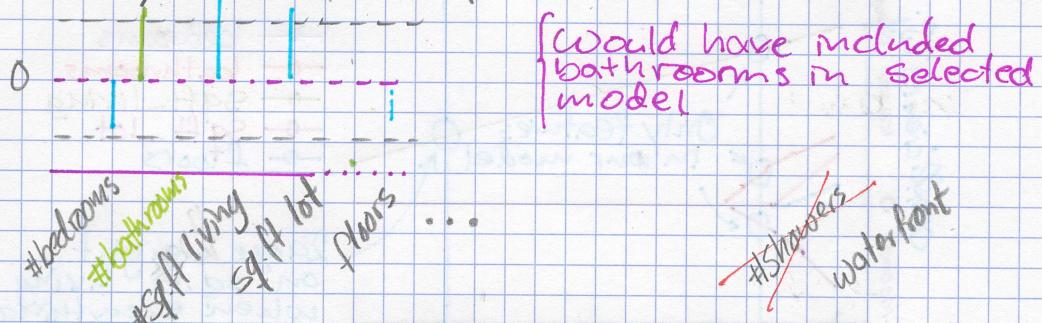
Why don't we just set small ridge coefficients to 0?



Let's look at two related features.. (in green)

Nothing measuring bathrooms was included!

If only one of the features had been included



Can regularization lead directly to sparsity?

Try this cost instead of ridge...

Total cost =

$$\underbrace{\text{measure of fit} + \lambda \text{ measure of magnitude of coefficients}}_{\text{RSS}(\mathbf{w})}$$

$$\|\mathbf{w}\|_1 = |\mathbf{w}_0| + \dots + |\mathbf{w}_D|$$

Lasso regression
(a.k.a. L₁ regularized regression)

Leads to sparse solutions!

Lasso regression:
L₁ regularized regression

Just like ridge regression solution is governed by a continuous parameter λ

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

λ tuning parameter = balance of fit and sparsity

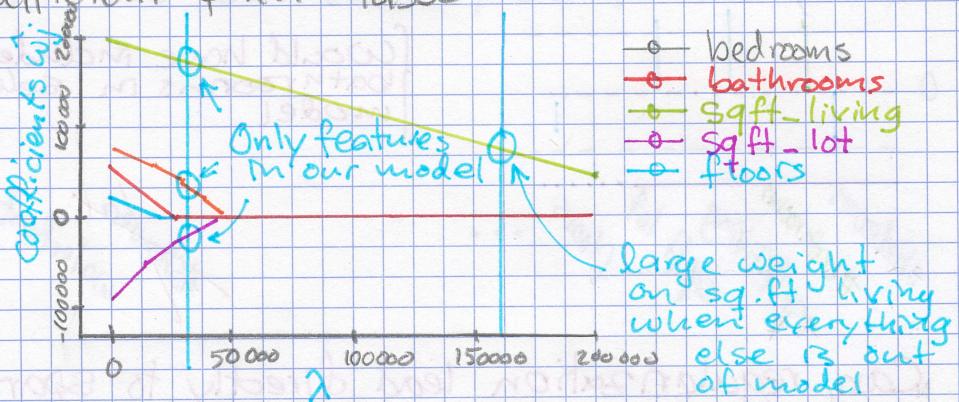
If $\lambda = 0$: $\hat{\mathbf{w}}^{\text{lasso}} = \hat{\mathbf{w}}^{\text{LS}}$ (unregularized solution)

If $\lambda = \infty$: $\hat{\mathbf{w}}^{\text{lasso}} = \mathbf{0}$

If γ in between

$$0 \leq \| \hat{w}^{\text{lasso}} \|_1 \leq \| \hat{w}^{\text{LS}} \|_1$$

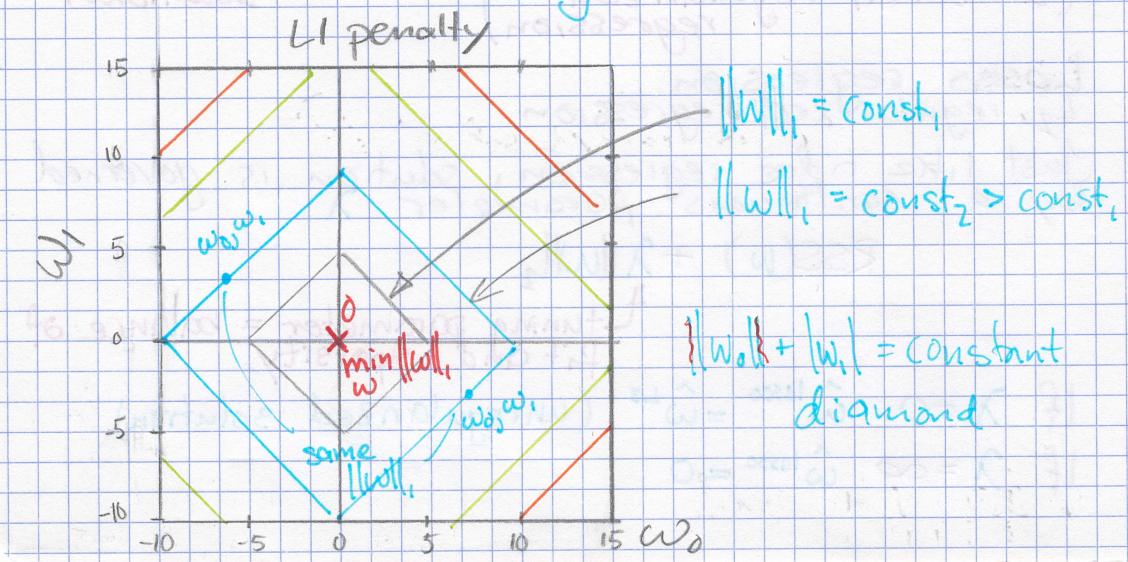
Coefficient path - lasso



$$RSS(\omega) + \lambda \|\omega\|_2^2 = \sum_{i=1}^N (y_i - \omega_0 h_0(x_i) - \omega_1 h_1(x_i))^2 + \lambda (\omega_0^2 + \omega_1^2)$$

$$\underbrace{\text{RSS}(\omega)}_{\text{ellipses}} + \lambda \underbrace{\| \omega \|_2^2}_{\text{circles}}$$

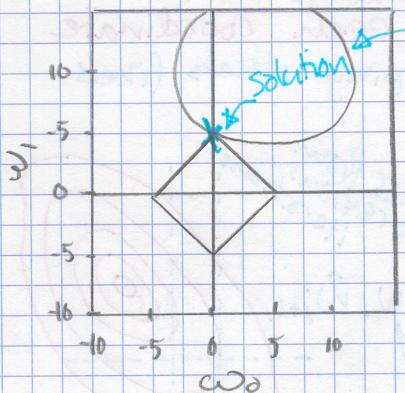
weighting



$$RSS(w) + \lambda \|w\|_1 = \sum_{i=1}^n (y_i - w_0 h_0(x_i) - w_1 h_1(x_i))^2 + \lambda (|w_0| + |w_1|)$$

RSS + $\lambda \|w\|_1$

ellipses diamonds



For a specific value of λ ,

Fitting the lasso regression model (for given λ value)

How we optimized past objectives
To solve for \hat{w} previously took gradient of total cost objective and either:

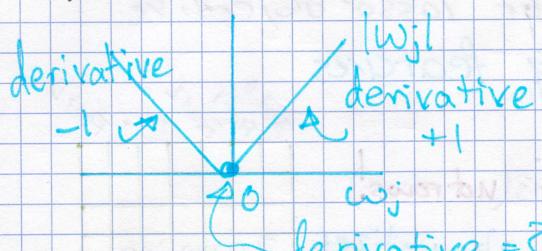
- 1) Derived closed form solution
- 2) Used in gradient descent algorithm

Optimizing the lasso objective

Lasso total cost $\text{RSS}(\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1$

Issues:

1) What's the derivative of $|y|$?



derivative =? no derivative exists

2) Even if we could compute derivative, no closed-form solution

Aside 1: Coordinate descent

Coordinate descent

Goal: Minimize some function g

$$\min_w g(w)$$

$$g(w) = g(w_0, w_1, \dots, w_D)$$

Often hard to find minimum for all coordinates, but easy for each coordinate when keeping others fixed

Coordinate descent

Initialize $\hat{w} = 0$

while not converged

pick a coordinate j

$$\hat{w}_j \leftarrow \min_w g(\hat{w}_0, \hat{w}_1, \dots, \hat{w}_{j-1}, w, \hat{w}_{j+1}, \dots, \hat{w}_D)$$

values from previous iterations

just min over j^{th} coordinate



Comments on coordinate descent

how do we pick next coordinate?

- At random ("random" or "stochastic" coordinate descent), round robin, ...

No stepsize to choose!

Super useful approach for many problems

- Converges to optimum in some cases
(e.g. strongly convex)

- Converges for lasso objective

Aside 2: Normalizing features

Normalizing features

Scale training columns (not rows!) as

$$h_j(x_k) = \frac{h_j(x_k)}{\sqrt{\sum_{i=1}^N h_j(x_i)^2}} \quad \text{Normalizer } Z_j$$

	Feat 0	Feat 1	Feat 2	...	Feat D
Training features					

Apply same training scale factors to test data

				Test features

$$h_j(x_k) = \frac{h_j(x_k)}{\sqrt{\sum h_j(x_i)^2}} \quad \text{Normalizer } Z_j$$

apply to test point

Summing over training points

Aside 3: Coordinate descent for unregularized regression (for normalized features)

Optimizing least squares objective one coordinate at a time

$$RSS(w) = \sum_{i=1}^N \left(y_i - \sum_{j=0}^D w_j h_j(x_i) \right)^2$$

Normalized features

Fix all coordinates $w_{\cdot j}$ and take partial w.r.t w_j
all w_k for $k \neq j$

$$\frac{\partial}{\partial w_j} RSS(w) = 2 \sum_{i=1}^N h_j(x_i) \left(y_i - \sum_{j=0}^D w_j h_j(x_i) \right)$$

1d optimization
coordinate by coordinate

$$= -2 \sum_{i=1}^N h_j(x_i) \left(y_i - \sum_{k \neq j} w_k h_k(x_i) - w_j h_j(x_i) \right)$$

$$= -2 \sum_{i=1}^N h_j(x_i) \left(y_i - \underbrace{\sum_{k \neq j} w_k h_k(x_i)}_{\hat{\varphi}_j} \right) + \cancel{2 \sum_{i=1}^N h_j(x_i)^2}$$

$\hat{\varphi}_j$

$+ 2w_j \left[\sum_{i=1}^N h_j(x_i)^2 \right]$

$$= -2\hat{\varphi}_j + 2w_j$$

by definition of normalized features
 $= 1$

Set partial = 0 and solve

$$\frac{\partial}{\partial w_j} \text{RSS}(w) = -2q_j + 2w_j = 0$$

$$w_j = q_j$$

Coordinate descent for least squares regression

Initialize $\hat{w} = 0$ (or smartly...)

while not converged:

for $j = 0, 1, \dots, D$

$$\text{compute } q_j = \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(\hat{w}_{-j}))$$

$$\text{set: } \hat{w}_j = q_j$$

residual without feature j

prediction without feature j

Coordinate descent for lasso (for normalized features)

Initialize $\hat{w} = 0$ (or smartly)

while not converged:

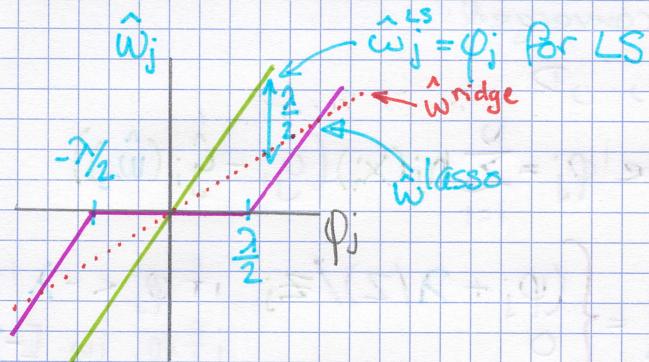
for $j = 0, 1, \dots, D$

$$\text{compute } q_j = \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(\hat{w}_{-j}))$$

$$\text{set } \hat{w}_j = \begin{cases} q_j + \lambda/2 & \text{if } q_j < -\lambda/2 \\ 0 & \text{if } q_j \in [-\lambda/2, \lambda/2] \\ q_j - \lambda/2 & \text{if } q_j > \lambda/2 \end{cases}$$

Soft thresholding

$$\hat{w}_j = \begin{cases} \varphi_j + \lambda/2 & \text{if } \varphi_j < -\lambda/2 \\ 0 & \text{if } \varphi_j \text{ in } [-\lambda/2, \lambda/2] \\ \varphi_j - \lambda/2 & \text{if } \varphi_j > \lambda/2 \end{cases}$$



Convergence criteria When to stop?

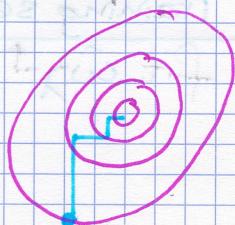
For convex problems, will start to take smaller and smaller steps

Measure size of steps taken in a full loop over all features

- Stop when max step < ε

Other lasso solvers

Classically: Least angle regression (LARS) [Efron et al. '04]



Then: Coordinate descent algorithm [Fu '98, Friedman, Hastie, & Tibshirani '08]

Now:

- Parallel CD (e.g. Shotgun, [Bradley et al. '11])
- Other parallel learning approaches for linear models
 - Parallel stochastic gradient descent (SGD) (e.g. Hogwild! [Niu et al. '11])
 - Parallel independent solutions then averaging [Zhang et al. '12]
- Alternating directions method of multipliers (ADMM) [Boyd et al. '11]

Coordinate descent for lasso (for unnormalized features)

Precompute $z_j = \sum_{i=1}^n h_j(x_i)^2$

Initialize $\hat{\omega} = 0$

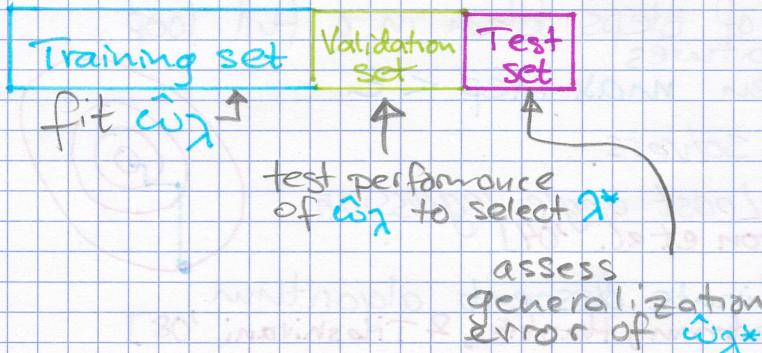
while not converged:

- for $j = 0, 1, \dots, D$

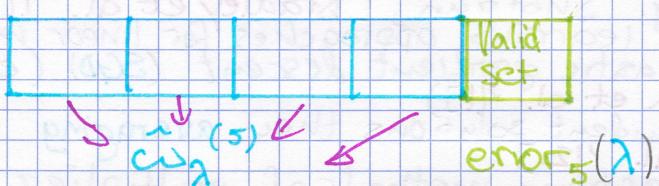
compute $\varphi_j = \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(\hat{\omega}))$

$$\text{set } \hat{\omega}_j = \begin{cases} (\varphi_j + \lambda/2)/z_j & \text{if } \varphi_j < -\lambda/2 \\ 0 & \text{if } \varphi_j \in [-\lambda/2, \lambda/2] \\ (\varphi_j - \lambda/2)/z_j & \text{if } \varphi_j > \lambda/2 \end{cases}$$

How to choose λ
If sufficient amount of data



k -fold cross validation



For $k = 1, \dots, K$

- Estimate $\hat{\omega}_\lambda^{(k)}$ on the training blocks
- Compute error on validation block $error_k(\lambda)$

Compute average error $CV(\lambda) = \frac{1}{K} \sum_{k=1}^K \text{error}_k(\lambda)$

Choosing λ via cross validation

Cross validation is choosing the λ that provides best predictive accuracy

Tends to favor less sparse solutions and thus smaller λ than optimal choice for feature selection

c.f. "Machine Learning: A Probabilistic Perspective", Murphy, 2012 for further discussion

Practical concerns with Lasso

Debiasing Lasso

Lasso shrinks coefficients relative to LS solution
→ more bias, less variance

- Can reduce bias as follows:
- 1. Run Lasso to select features
- 2. Run least squares regression with only selected features

"Relevant" features no longer shrunk relative to LS fit of same reduced model

Issues with standard lasso objective

1. With group of highly correlated features, lasso tends to select amongst them arbitrarily
- often prefers to select all together
2. Often empirically ridge has better predictive performance than lasso, but lasso leads to sparser solution

Elastic net aims to address these issues
- hybrid between lasso and ridge regression
- uses L_1 and L_2 penalties

See Zou & Hastie '05 for further discussion

Impact of feature selection and lasso

Lasso has changed machine learning, statistics & electrical engineering

But, for feature selection in general, be **Careful about interpreting selected features**

- Selection only considers features included
- sensitive to correlations between features
- result depends on algorithm used
- there are theoretical guarantees for lasso under certain conditions