

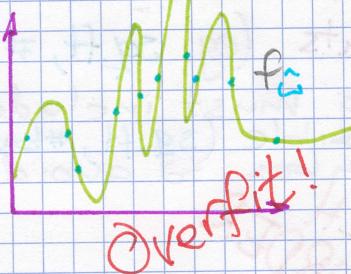
## Ridge Regression

Regulating overfitting when using many features

## Overfitting of polynomial regression

Flexibility of high-order polynomials

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$



## Symptom of overfitting

Often overfitting associated with very large estimated parameters  $w$

## Overfitting of linear regression models more generically

Overfitting with many features

Not unique to polynomial regression, but also if lots of inputs ( $d$  large)

Or generically lots of features ( $D$  large)

$$y_i = \sum_{j=0}^D w_j h_j(x_i) + \epsilon_i$$

- Square feet
- # bathrooms
- # bedrooms
- Lot size
- Year built
- ...

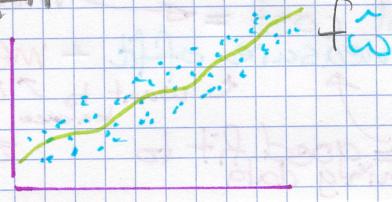
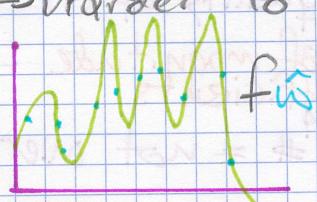
How does # of observations influence overfitting?

Few observations ( $N$  small)

→ rapidly overfit as model complexity increases

Many observations ( $N$  very large)

→ harder to overfit



1 input (e.g., Sq. ft.)

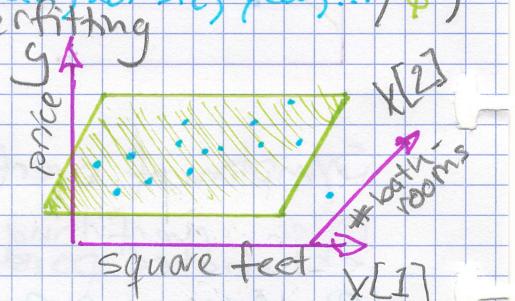
Data must include representative examples of all possible (sq.ft., \$) pairs to avoid overfitting

HARD

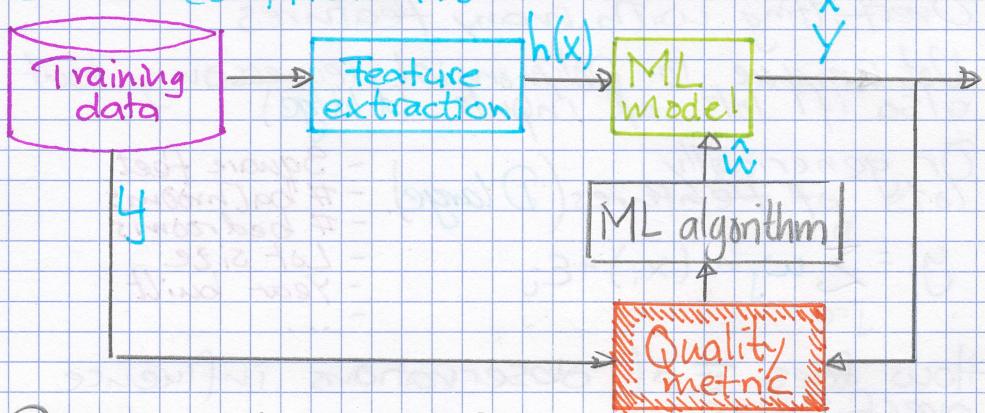
d inputs (e.g., sq.ft., #bath, #bed, lot size, year, ...)

Data must include examples of all possible (sq.ft., #bath, #bed, lot size, year, ..., \$) combos to avoid overfitting

MUCH HARDER



Adding term to cost-of-fit to prefer small coefficients



Desired total cost format

want to balance

- i. How well function fits data
- ii. Magnitude of coefficients

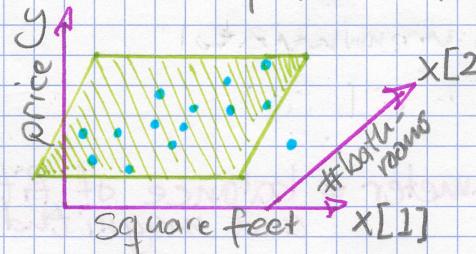
Total cost = measure quality of fit

$$\text{measure of fit} + \text{measure of magnitude}$$

↑ ↑ want to balance of coefficients

Small # = good fit to training data      Small # = not overfit

Measure of fit to training data



$$RSS(\mathbf{w}) = \sum_{i=1}^N (y_i - h(x_i)^T \mathbf{w})^2$$

$$= \sum_{i=1}^N (y_i - \hat{y}_i(\mathbf{w}))^2$$

pred. value using  
 $\mathbf{w}$

small RSS  $\rightarrow$  model fitting  
training data well

Measure of the magnitude  
of regression coefficient

What summary # is indicative of size of  
regression coefficients?

- Sum?  $w_0 = 1,527,301, w_1 = -1,605,253$   
 $w_0 + w_1 = \text{small } \#$

- Sum of absolute value?

$$|w_0| + |w_1| + \dots + |w_d| = \sum_{j=0}^d |w_j| \leq \|w\|_1, \quad L_1 \text{ norm}$$

- Sum of squares ( $L_2$  norm)

$$w_0^2 + w_1^2 + \dots + w_d^2 = \sum_{j=0}^d w_j^2 \leq \|w\|_2^2, \quad L_2 \text{ norm}$$

Total cost =

$$\underbrace{RSS(\mathbf{w})}_{\text{measure of fit}} + \underbrace{\|w\|_2^2}_{\text{measure of magnitude of coefficients}}$$

Consider resulting objective

What if  $\hat{w}$  selected to minimize

$$RSS(w) + \lambda \|w\|_2^2$$

↑  
tuning parameter = balance of fit  
and magnitude

If  $\lambda = 0$ :

reduces to minimizing  $RSS(w)$ , as before  
(old solution)  $\rightarrow \hat{w}^{LS}$   
least squares

If  $\lambda = \infty$ :

For solutions where  $\hat{w} \neq 0$ , then total cost is  $\infty$   
If  $w = 0$ , then total cost =  $RSS(0)$   
 $\rightarrow$  solution is  $\hat{w} = 0$

If  $\lambda$  in between  
Then  $0 \leq \|w\|_2^2 \leq \|\hat{w}^{LS}\|_2^2$

Ridge regression  
(a.k.a L<sub>2</sub> regularization)

Bias-variance tradeoff

Large  $\lambda$ :  
high bias, low variance  
(eg,  $\hat{w} = 0$  for  $\lambda = \infty$ )

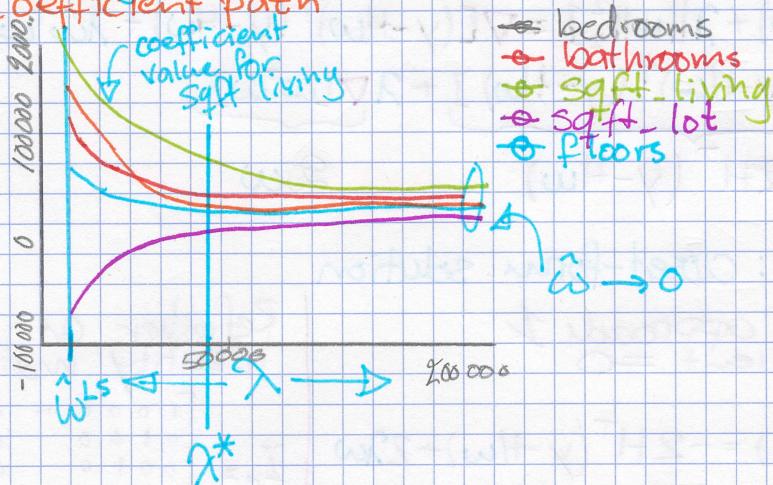
{ In essence,  $\lambda$   
controls model complexity

Small  $\lambda$ :

low bias, high variance  
(eg, standard least squares (RSS) fit of  
high-order polynomial for  $\lambda = 0$ )

## The ridge coefficient path

### Coefficient path



Fitting the ridge regression model (for given  $\lambda$  value)

Step 1: Rewrite total cost in matrix notation

Rewrite magnitude of coefficients in vector notation

$$\|w\|_2^2 = w_0^2 + w_1^2 + w_2^2 + \dots + w_D^2$$

$$= \begin{matrix} \boxed{\phantom{0}}, \boxed{\phantom{0}}, \boxed{\phantom{0}}, \dots, \boxed{\phantom{0}} \\ w_0, w_1, w_2, \dots, w_D \end{matrix} \begin{matrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{matrix}$$

$$= w^T w$$

Putting it all together In matrix form, ridge regression cost is:

$$RSS(w) + \lambda \|w\|_2^2$$

$$= (y - Hw)^T (y - Hw) + \lambda w^T w$$

Gradient of ridge regression cost

$$\nabla [RSS(\omega) + \lambda \|\omega\|_2^2] = \nabla [(y - H\omega)^T (y - H\omega) + 2\lambda \omega^T \omega]$$

$$= \underbrace{\nabla [(y - H\omega)^T (y - H\omega)]}_{-2H^T(y - H\omega)} + \underbrace{2\lambda \nabla [\omega^T \omega]}_{2\lambda \omega}$$

Approach 1: closed-form solution

Step 3: approach 1  
set gradient  $\equiv 0$

$$\nabla \text{cost}(\omega) = -2H^T(y - H\omega) + 2\lambda \omega$$

$$= -2H^T(y - H\omega) + 2\lambda I\omega$$

Refresher on identity matrix

$$I_n = \begin{bmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 & & 0 \\ 0 & 0 & 1 & 0 & & 0 \\ 0 & 0 & 0 & 1 & & \vdots \\ \vdots & & & & & 1 \\ 0 & 0 & 0 & 0 & & 1 \end{bmatrix}$$

Ridge closed-form solution:

$$\nabla \text{cost}(\omega) = -2H^T(y - H\omega) + 2\lambda I\omega = 0$$

$$-H^T y + H^T H\omega + \lambda I\hat{\omega} = 0.$$

$$H^T H\hat{\omega} + \lambda I\hat{\omega} = H^T y$$

$$(H^T H + \lambda I)\hat{\omega} = H^T y$$

$$\hat{\omega} = (H^T H + \lambda I)^{-1} H^T y$$

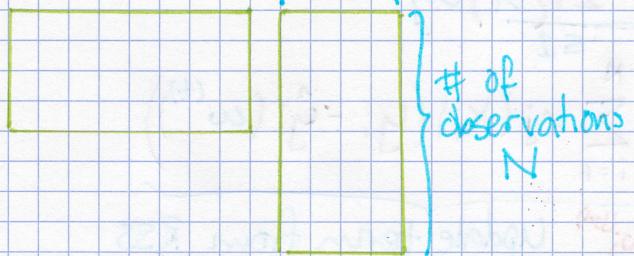
Interpreting ridge closed-form solution

$$\hat{\omega}^{\text{ridge}} = (H^T H + \lambda I)^{-1} H^T y$$

If  $\lambda = 0$ :  $\hat{\omega}^{\text{ridge}} = (H^T H)^{-1} H^T y = \hat{\omega}^{\text{LS}}$   $\leftarrow$  old solution!

If  $\lambda = \infty$ :  $\hat{\omega}^{\text{ridge}} = 0$   $\uparrow$  because it's like dividing by  $\infty$

$$\hat{w}^{LS} = \underbrace{(H^T H)^{-1}}_{\# \text{ of features } D} H^T y$$



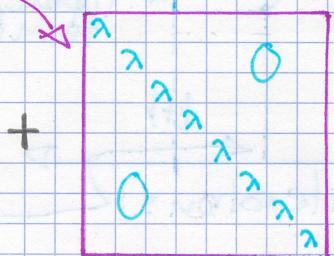
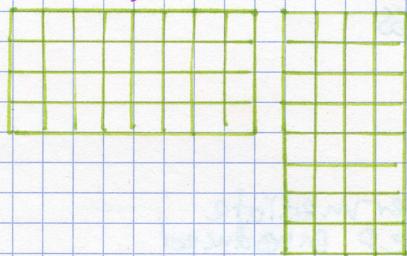
Invertible if  
In general,  
(# linearly  
independent obs.)  
 $N > D$

Complexity of  
inverse  
 $\mathcal{O}(D^3)$

really important for  
large  $D$  (lots of features)

Invertible if  
Always if  $\lambda > 0$ ,  
even if  $N < D$

$$\hat{w} = \underbrace{(H^T H + \lambda I)^{-1}}_{\# \text{ of features } D} H^T y$$



Complexity  
of inverse  
 $\mathcal{O}(D^3)$   
big for large  
 $D$ !

" $\lambda I$  is making  $H^T H + \lambda I$  more  
regular"  $\rightarrow$  "regularized"

## Approach 2: gradient descent

### Step 3: approach 2: Gradient descent

Elementwise ridge regression gradient descent  
algorithm

$$\nabla \text{cost}(w) = -2H^T(y - Hw) + 2\lambda w$$

Update to  $j$ th feature weight:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta * \left[ -2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)})) + 2\lambda w_j^{(t)} \right]$$

new term, comes from the  
 $j$ th component of  $\nabla \text{cost}(w)$

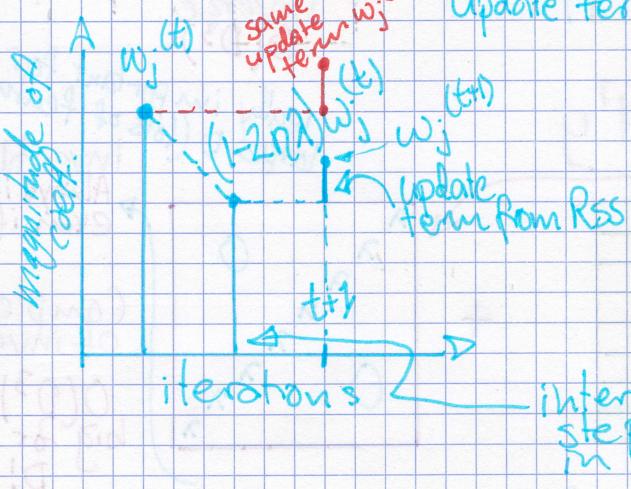
Equivalently:  $\frac{2\eta\lambda}{1-2\eta\lambda} \leq 1$

$$w_j^{(t+1)} \leftarrow (1 - \frac{2\eta\lambda}{1-2\eta\lambda}) w_j^{(t)}$$

$$+ 2\eta \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$$

Update term from RSS

Previously, just  
RSS (LS)



intermediate  
step introduced  
in ridge regression

Summary of ridge regression algorithm

init  $w^{(1)} = 0$  (or randomly, or smartly),  $t=1$

while  $\|\nabla \text{RSS}(w^{(t)})\| > \epsilon$

for  $j = 0, \dots, D$

$$\text{partial}[j] = -2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$$

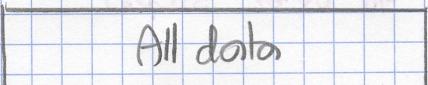
$$w_j^{(t+1)} \leftarrow (1 - 2\eta\lambda) w_j^{(t)} - \eta \text{partial}[j]$$

$t \leftarrow t + 1$

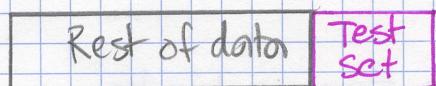
How to choose  $\alpha$   
if sufficient amount of data



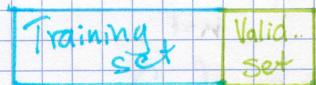
Start with smallish dataset



Still form test set and hold out



use for both training  
and validation, but not so  
naively



Is validation set enough to compare  
performance of  $\hat{w}_\alpha$  across  $\alpha$  values?

NO

Choosing the validation set



Which subset should I use?

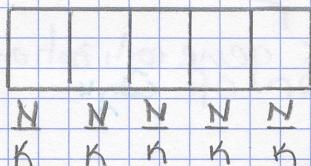
ALL! average performance  
(over all choices)

## K-fold cross validation

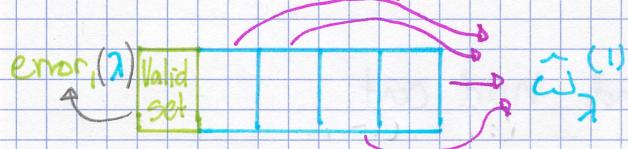
Rest of data

### Preprocessing:

Randomly assign data to K groups



(use same split of data for all other steps)



For  $k=1, \dots, K$

1. Estimate  $\hat{w}_j^{(k)}$  on the training blocks
2. Compute error on validation block:  $\text{error}_k(\lambda)$

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

Repeat procedure for each choice of  $\lambda$

Choose  $\lambda^*$  to minimize  $CV(\lambda)$

What value of  $K$ ?

Formally the best approximation occurs for validation sets of size 1 ( $K=N$ )

leave-one-out  
cross validation

Computationally intensive

- requires computing  $N$  fits of model per  $\lambda$

Typically,  $K=5$  or  $10$

5-fold CV

10-fold CV

## How to handle the intercept

Recall multiple regression model

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \epsilon_i$$
$$= \sum_{j=0}^D w_j h_j(x_i) + \epsilon_i$$

feature 1 =  $h_0(x)$  ... often 1 (constant)

feature 2 =  $h_1(x)$  ... e.g.,  $x[1]$

feature 3 =  $h_2(x)$  ... e.g.,  $x[2]$

feature D+1 =  $h_D(x)$  ... e.g.,  $x[d]$

In matrix notation for N observations

$$\mathbf{y} = \mathbf{H}\mathbf{w} + \boldsymbol{\epsilon}$$
$$\begin{array}{|c|c|c|} \hline & \begin{array}{|c|c|c|c|c|} \hline & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & \dots & 1 \\ \hline & 1 & 1 & \dots & 1 \\ \hline \end{array} & = & \begin{array}{|c|} \hline w_0 \\ \hline \vdots \\ \hline w_d \\ \hline \end{array} + \begin{array}{|c|} \hline \boldsymbol{\epsilon} \\ \hline \vdots \\ \hline \end{array} \\ \hline \end{array}$$

Do we penalize intercept?

Standard ridge regression cost

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↑ strength of penalty

Encourages intercept  $w_0$  to also be small

Do we want a small intercept?  
(conceptually) not indicative of overfitting..

Option 1: Don't penalize intercept

Modified ridge regression cost:

$$\text{RSS}(w_0, w_{\text{rest}}) + \lambda \|w_{\text{rest}}\|_2^2$$

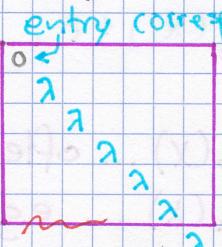
How to implement this in practice?

Option 1: Don't penalize intercept

- Closed-form solution -

$$\hat{\omega} = (H^T H + \lambda I)^{-1} H^T y$$

entry corresponding to  $w_0$  index



Option 1: Don't penalize intercept  
Gradient descent algorithm

while  $\|\nabla \text{RSS}(w)\| > \epsilon$

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

$$\text{partial}[j] = -2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$$

if  $j == 0$ :

$$w_0^{(t+1)} \leftarrow w_0^{(t)} - \eta \text{partial}[j] \quad \begin{matrix} \leftarrow \text{old LS} \\ \text{update} \\ \text{no shrinkage} \end{matrix}$$

else for all other features

$$w_j^{(t+1)} \leftarrow (1 - 2\eta)w_j^{(t)} - \eta \text{partial}[j]$$

$t \leftarrow t + 1$

ridge regression

Option 2: Center data first

If data are first centered about 0, then favoring small intercept not so worrisome

Step 1: Transform  $y$  to have 0 mean

Step 2: Run ridge regression as normal  
(closed-form or gradient algorithm)