

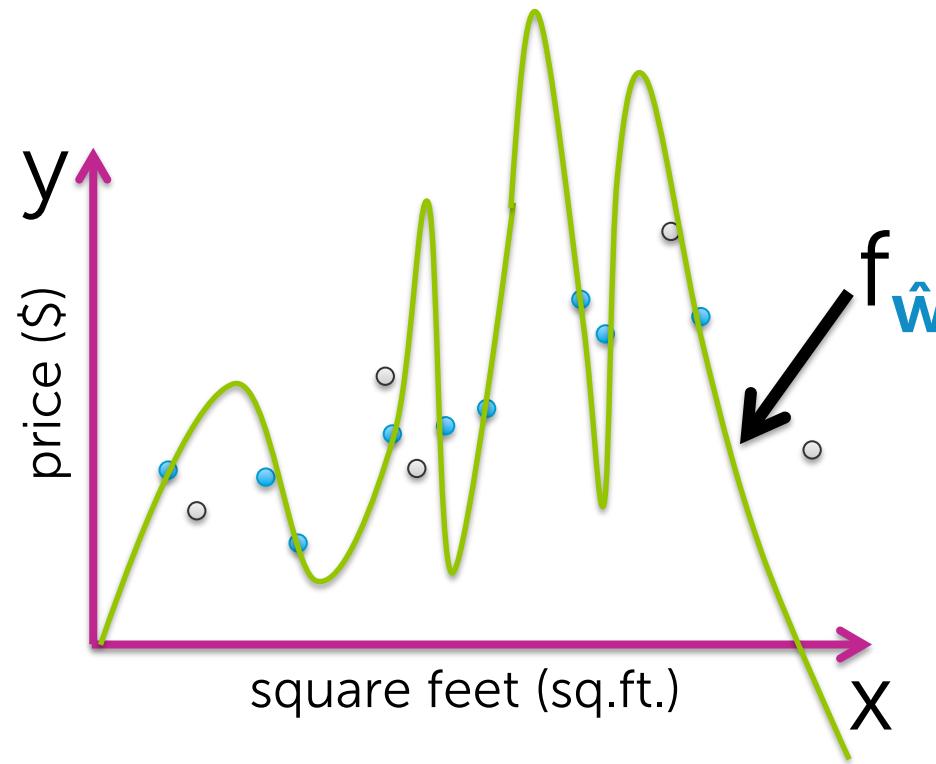
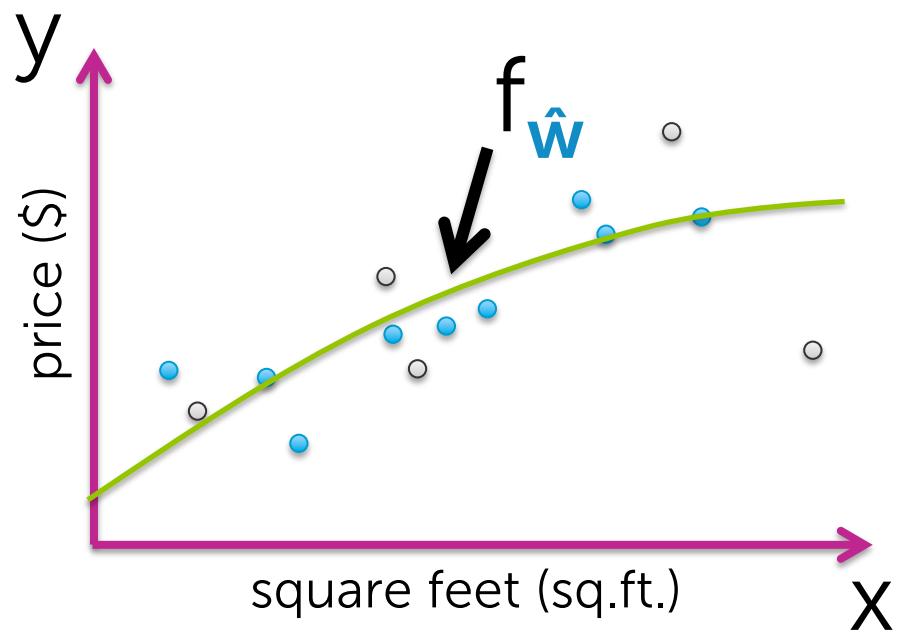
Ridge Regression:

Regulating overfitting when using many features

Overfitting of polynomial regression

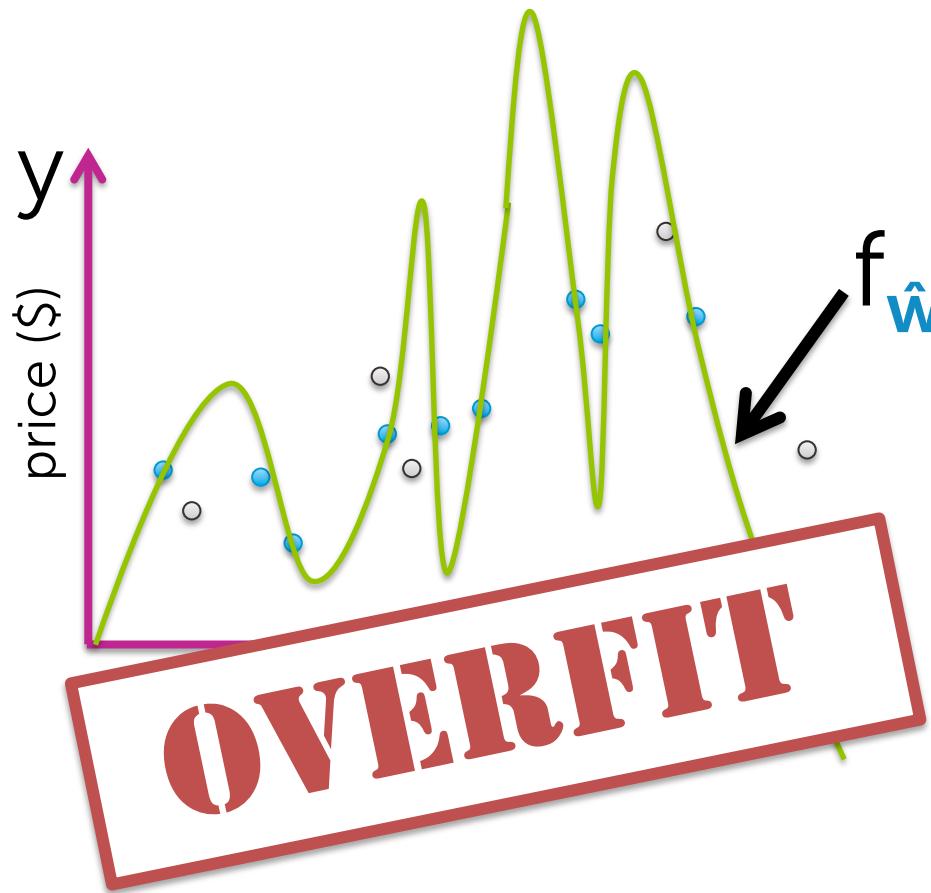
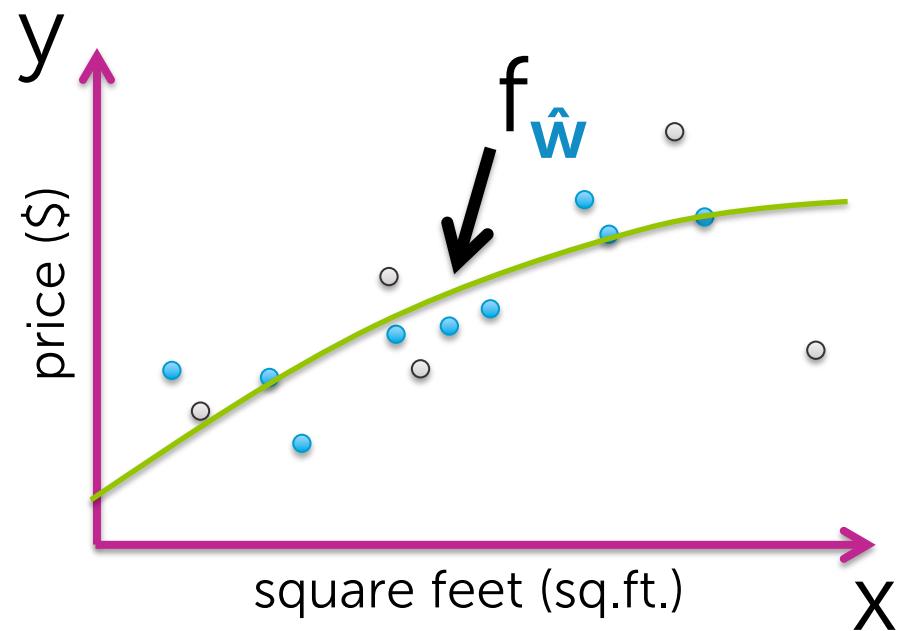
Flexibility of high-order polynomials

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



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 \hat{w}

OVERFITTING

DEMO

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)

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

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

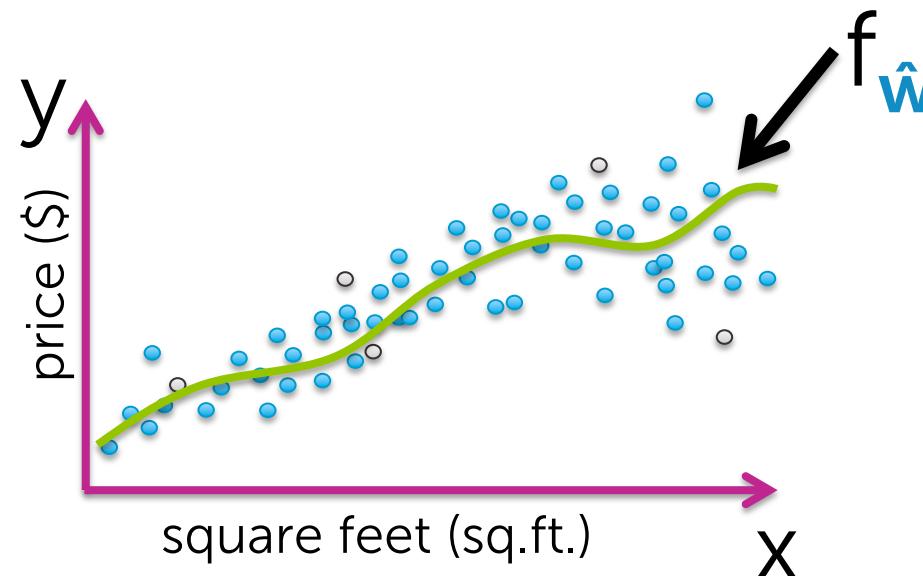
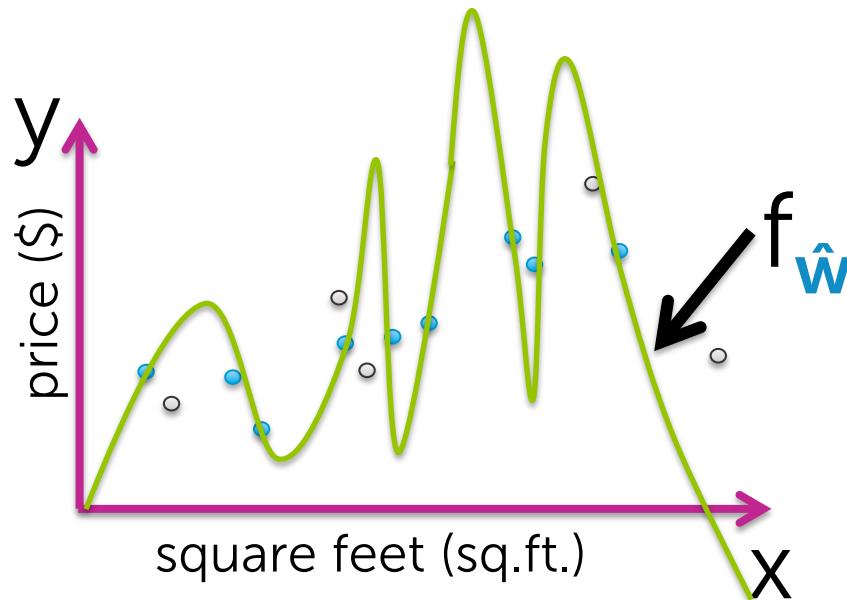
How does # of observations influence overfitting?

Few observations (N small)

→ rapidly overfit as model complexity increases

Many observations (N very large)

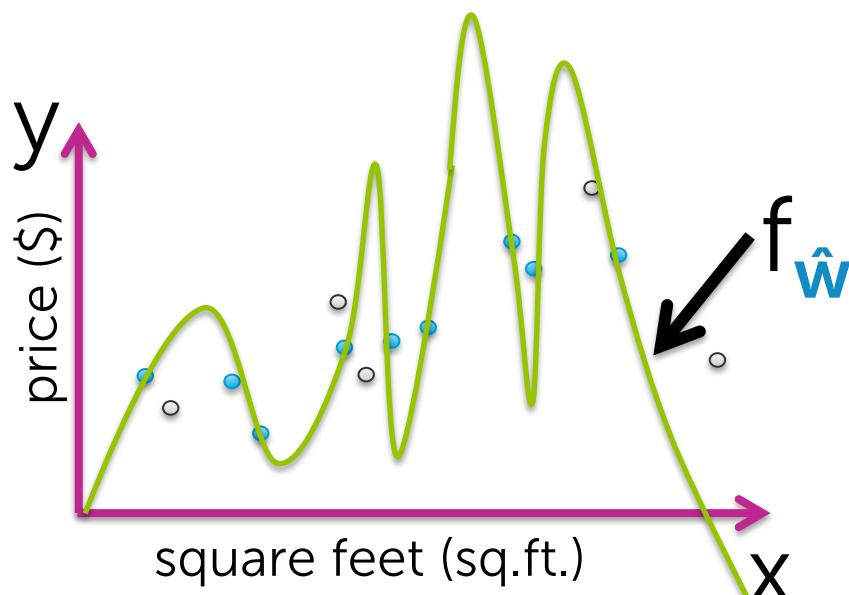
→ harder to overfit



How does # of inputs influence overfitting?

1 input (e.g., sq.ft.):

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

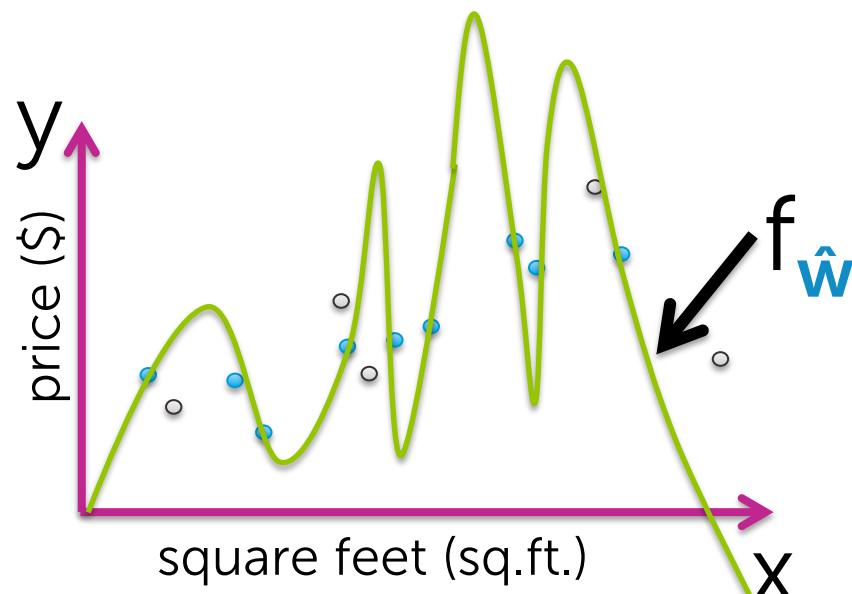


How does # of inputs influence overfitting?

1 input (e.g., sq.ft.):

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

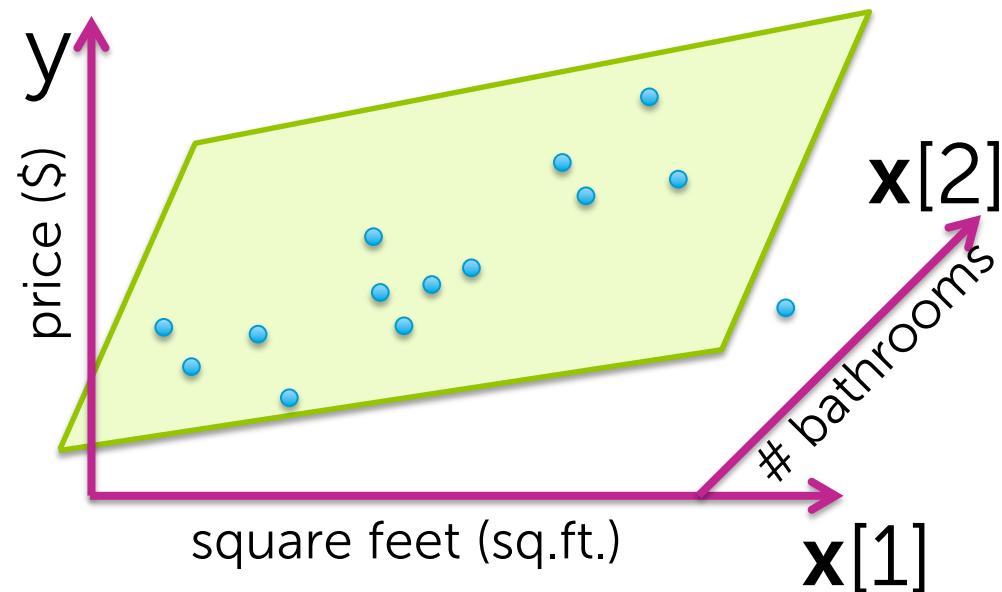
HARD



How does # of inputs influence overfitting?

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

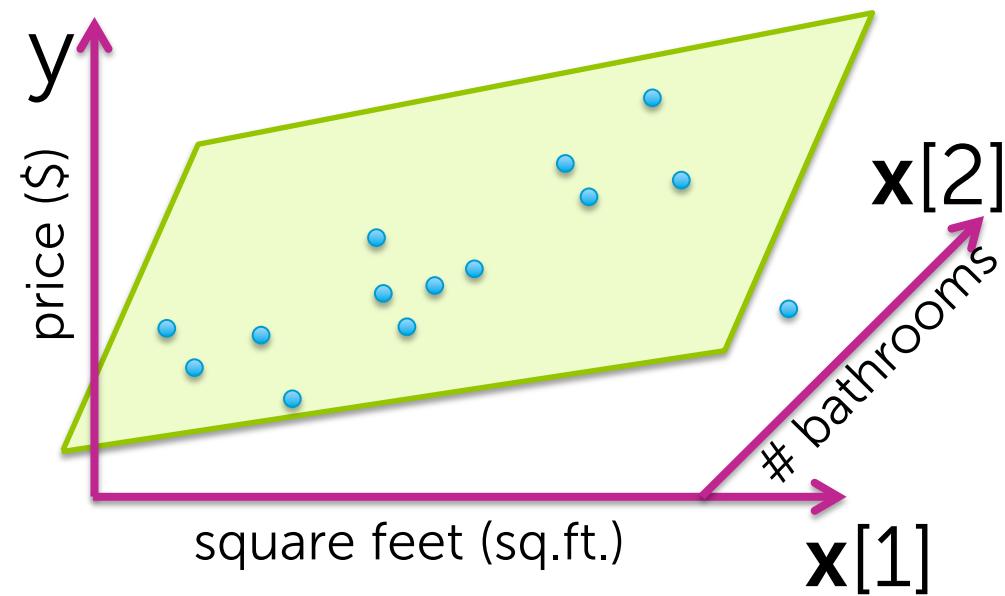


How does # of inputs influence overfitting?

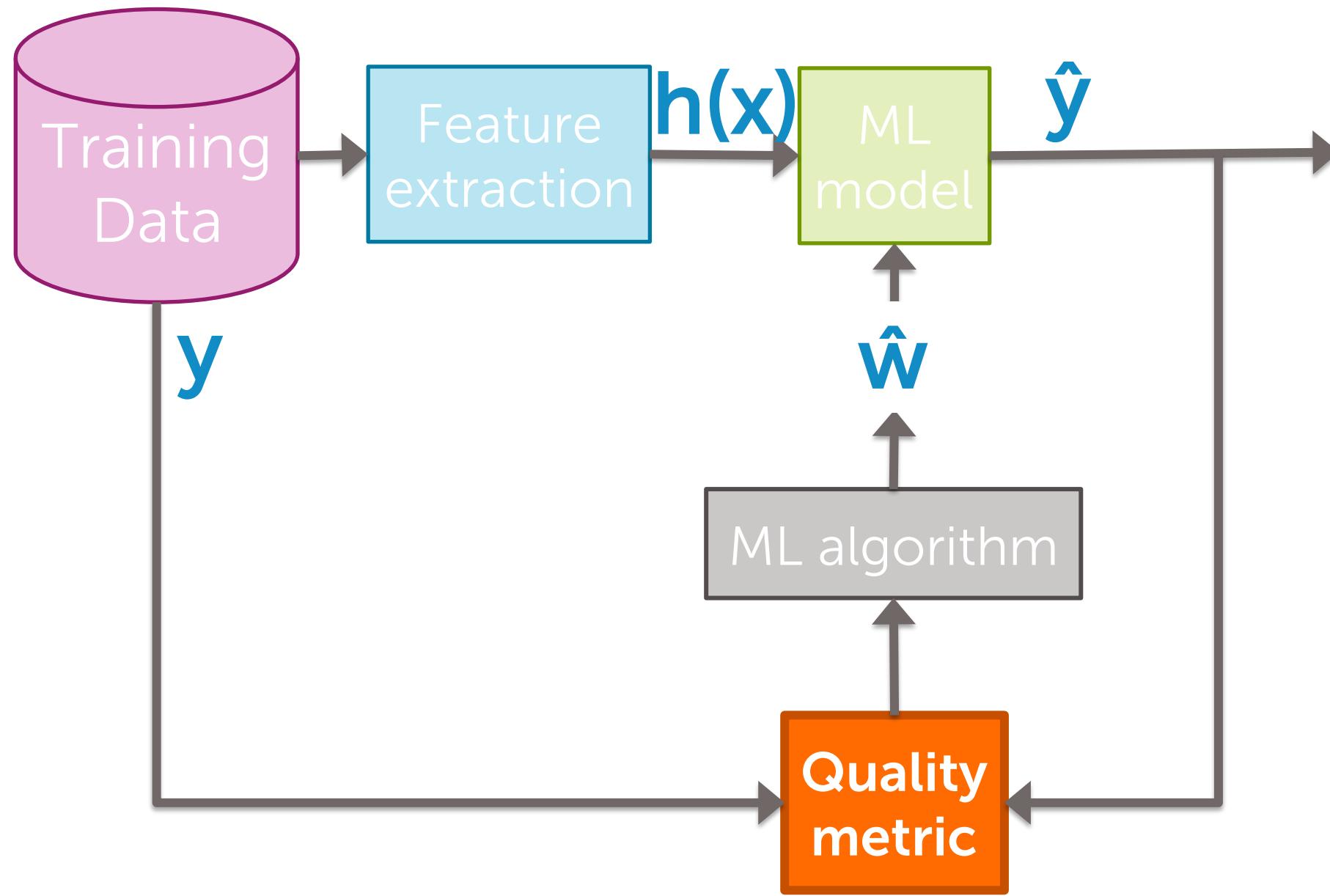
d inputs (e.g., sq.ft., #bath, #bed, lot size, year,...):

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

MUCH!!!
HARDER



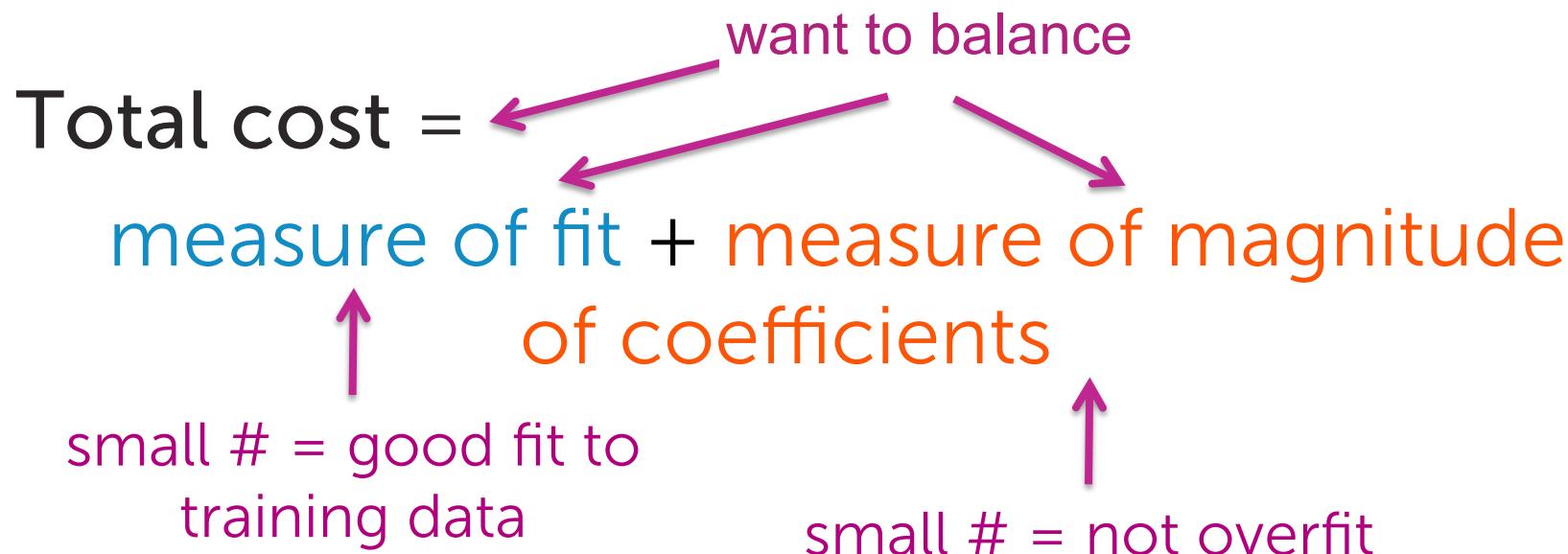
Adding term to cost-of-fit
to prefer small coefficients



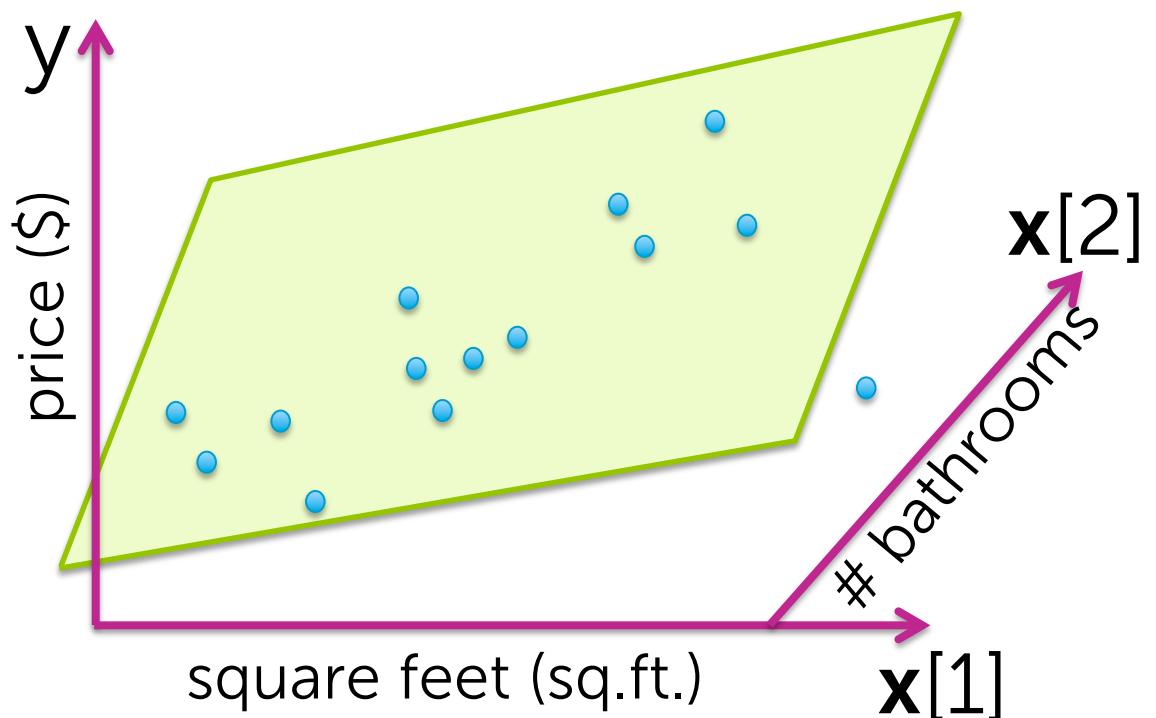
Desired total cost format

Want to balance:

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



Measure of fit to training data



$$\begin{aligned} \text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2 \\ &= \sum_{i=1}^N (y_i - \hat{y}_i(\mathbf{w}))^2 \end{aligned}$$

pred. value using w

small RSS \rightarrow model fitting training data well

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

- Sum?

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

- Sum? $w_0 = 1,527,301 \quad w_1 = -1,605,253$
 $w_0 + w_1 = \text{small } \pm$
- Sum of absolute value?

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_p| = \sum_{j=0}^p |w_j| \triangleq \|w\|_1$

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_p| = \sum_{j=0}^p |w_j| \triangleq \|w\|_1, \quad L_1 \text{ norm}$

... next class

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_p| = \sum_{j=0}^p |w_j| \triangleq \|w\|_1, \quad L_1 \text{ norm}$

- Sum of squares (L_2 norm)

... next class

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_D| = \sum_{j=0}^D |w_j| \triangleq \|w\|_1, \quad L_1 \text{ norm}$

- Sum of squares (L_2 norm)
 $w_0^2 + w_1^2 + \dots + w_D^2 = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2$

... next class

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

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

- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_D| = \sum_{j=0}^D |w_j| \triangleq \|w\|_1, \quad L_1 \text{ norm}$

... next class

- Sum of squares (L_2 norm)
 $w_0^2 + w_1^2 + \dots + w_D^2 = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2 \quad L_2 \text{ norm}$

... this class

Consider specific total cost

Total cost =

measure of fit + measure of magnitude
of coefficients

Consider specific total cost

Total cost =

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



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

Consider resulting objective

What if \hat{w} selected to minimize

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

 tuning parameter = balance of fit and magnitude

If $\lambda=0$:

Consider resulting objective

What if \hat{w} selected to minimize

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda = 0$:

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

Consider resulting objective

What if \hat{w} selected to minimize

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda=0$:

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

If $\lambda=\infty$:

Consider resulting objective

What if \hat{w} selected to minimize

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda=0$:

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

If $\lambda=\infty$:

For solutions where $\hat{\mathbf{w}} \neq 0$, then total cost is ∞

If $\hat{\mathbf{w}}=0$, then total cost = $\text{RSS}(0)$ \rightarrow solution is $\hat{\mathbf{w}}=0$

Consider resulting objective

What if \hat{w} selected to minimize

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda=0$:

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

If $\lambda=\infty$:

For solutions where $\hat{\mathbf{w}} \neq 0$, then total cost is ∞

If $\hat{\mathbf{w}}=0$, then total cost = $\text{RSS}(0)$ \rightarrow solution is $\hat{\mathbf{w}}=0$

If λ in between:

Consider resulting objective

What if \hat{w} selected to minimize

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda=0$:

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

If $\lambda=\infty$:

For solutions where $\hat{w} \neq 0$, then total cost is ∞

If $\hat{w}=0$, then total cost = $\text{RSS}(0) \rightarrow$ solution is $\hat{w}=0$

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

Consider resulting objective

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

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

 tuning parameter = balance of fit and magnitude

Ridge regression
(a.k.a L_2 regularization)

Bias-variance tradeoff

Large λ :

high bias, low variance

(e.g., $\hat{w} = 0$ for $\lambda = \infty$)

Bias-variance tradeoff

Large λ :

high bias, low variance

(e.g., $\hat{w} = 0$ for $\lambda = \infty$)

Small λ :

low bias, high variance

(e.g., standard least squares (RSS) fit of
high-order polynomial for $\lambda = 0$)

Bias-variance tradeoff

Large λ :

high bias, low variance

(e.g., $\hat{w} = 0$ for $\lambda = \infty$)

In essence, λ
controls model
complexity

Small λ :

low bias, high variance

(e.g., standard least squares (RSS) fit of
high-order polynomial for $\lambda = 0$)

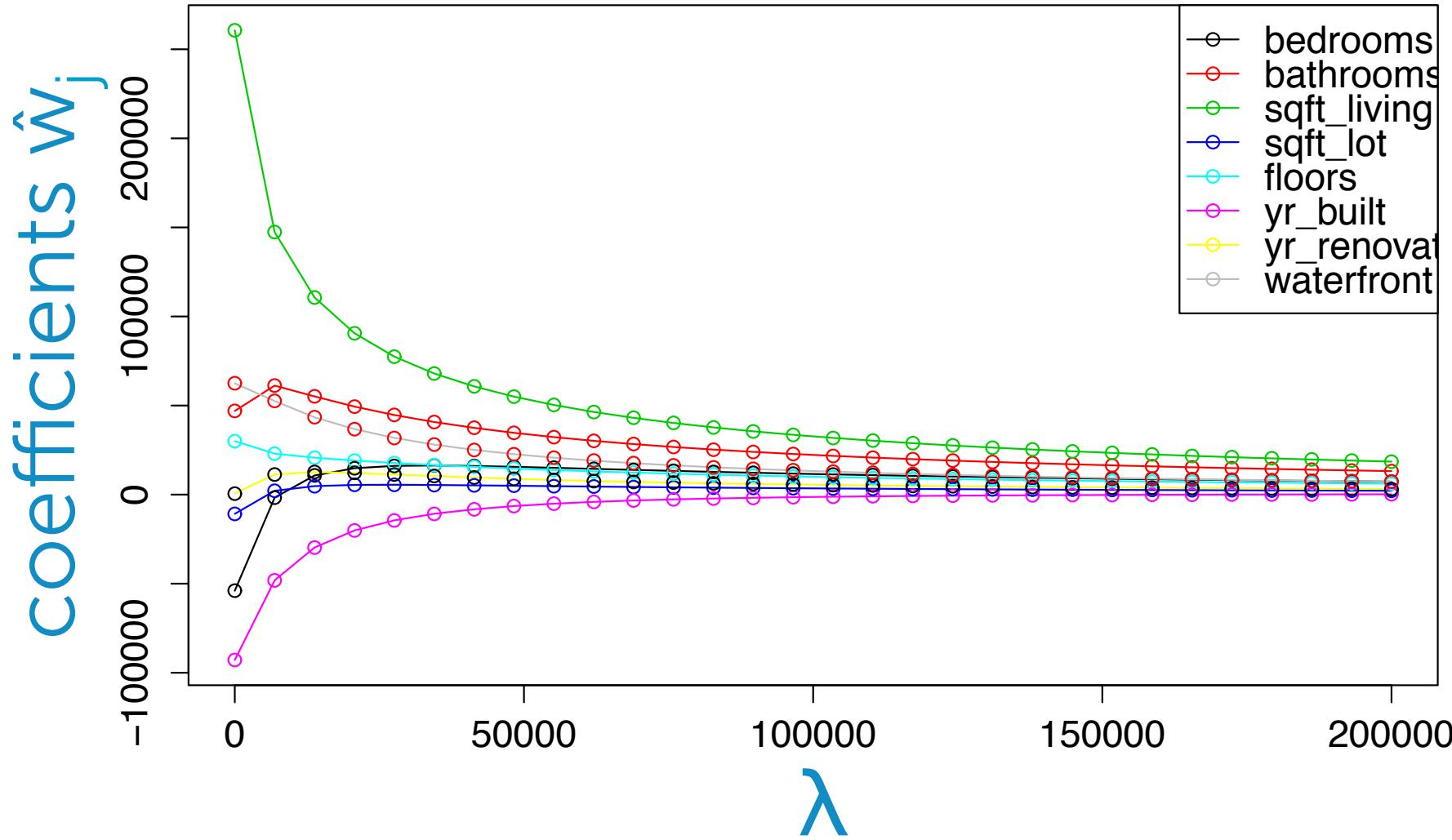
Revisit polynomial fit demo

What happens if we refit our high-order polynomial, but now using ridge regression?

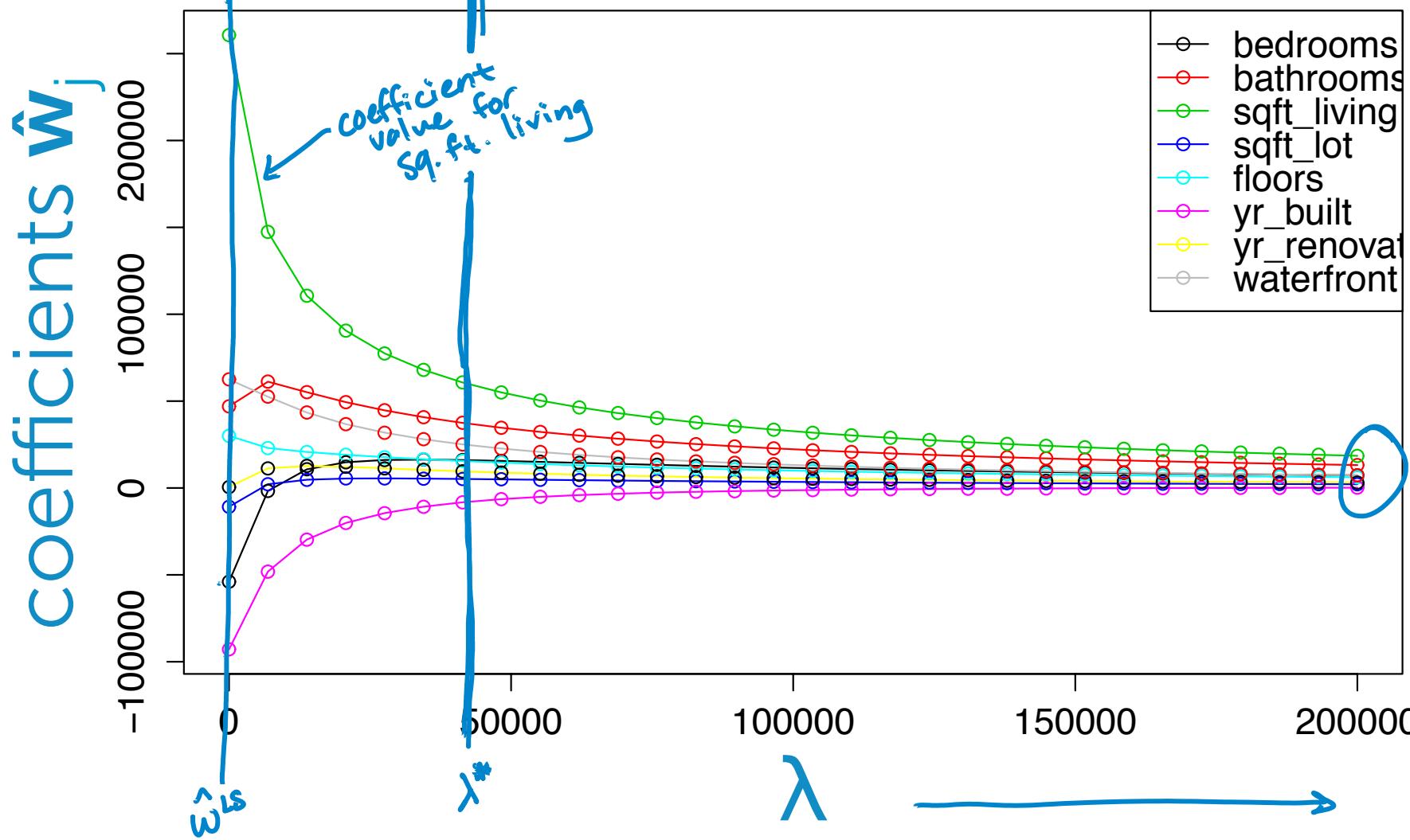
Will consider a few settings of λ ...

Back to notebook...

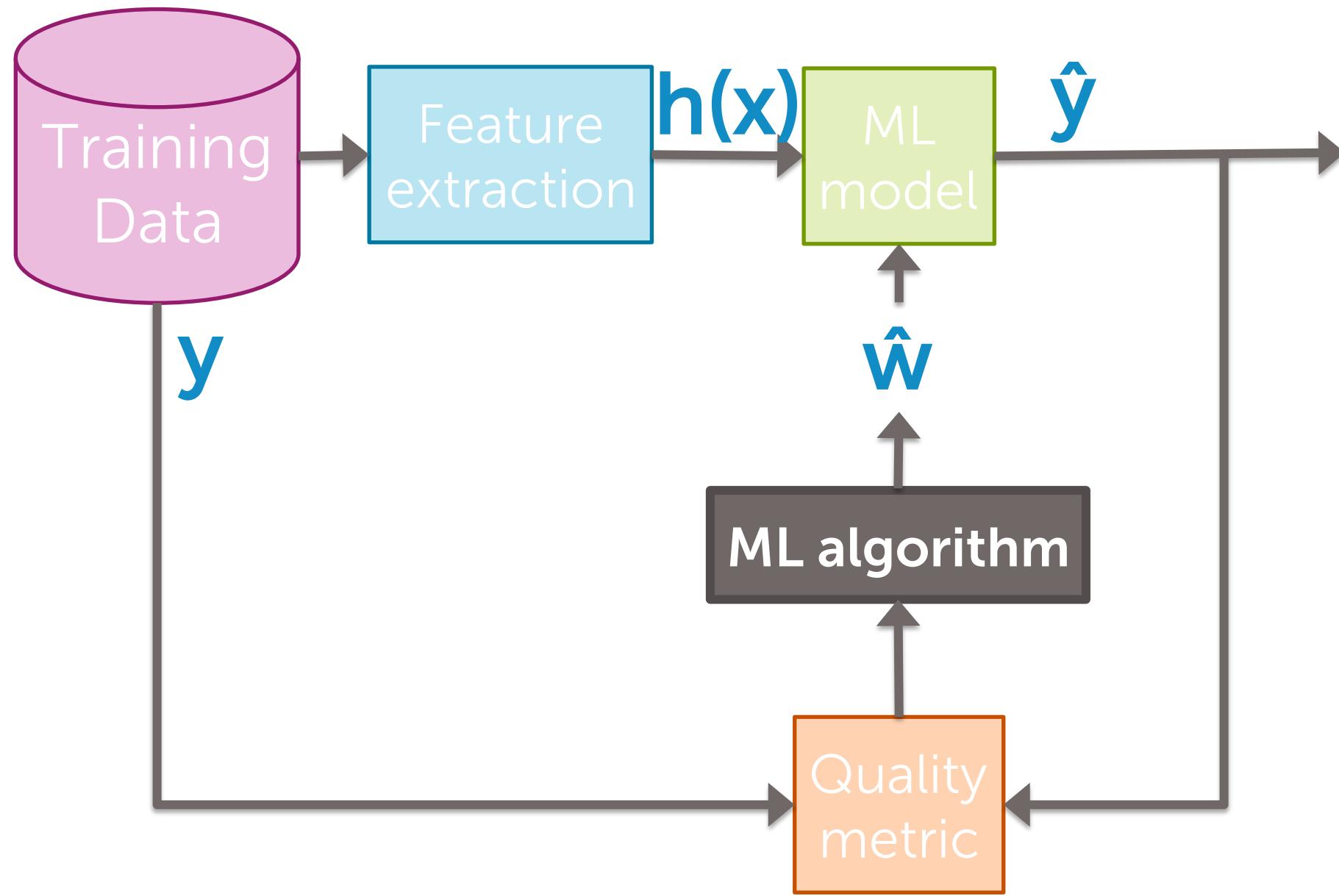
Coefficient path



Coefficient path



Fitting the ridge regression model
(for given λ value)



Step 1:
Rewrite total cost in matrix notation

Recall matrix form of RSS

Model for all N observations together

$$\mathbf{y} = \mathbf{H} \mathbf{w} + \boldsymbol{\epsilon}$$

The equation illustrates the matrix form of the linear regression model. On the left, the observed data vector \mathbf{y} is shown as a vertical column of pink squares. An equals sign follows. To the right of the equals sign is the matrix \mathbf{H} , which is a green square matrix with a black outline. To the right of \mathbf{H} is the weight vector \mathbf{w} , represented by a vertical column of blue squares. A plus sign follows. To the right of the plus sign is the error term $\boldsymbol{\epsilon}$, represented by a vertical column of grey squares.

Recall matrix form of RSS

$$\begin{aligned}\text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2 \\ &= (\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

Rewrite magnitude of coefficients in vector notation

$$\begin{aligned}\|w\|_2^2 &= w_0^2 + w_1^2 + w_2^2 + \dots + w_D^2 \\&= \begin{array}{c|c|c|c|c|c|c|c|c} \textcolor{blue}{w_0} & \textcolor{blue}{w_1} & \textcolor{blue}{w_2} & \dots & \textcolor{blue}{w_D} \end{array} \begin{array}{c} w_0 \\ w_1 \\ w_2 \\ \vdots \\ \vdots \\ w_D \end{array} \\&= w^T w\end{aligned}$$

Putting it all together

In matrix form, ridge regression cost is:

$$\begin{aligned} \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \\ = (\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \end{aligned}$$

Step 2: Compute the gradient

Gradient of ridge regression cost

$$\begin{aligned}\nabla [\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2] &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}] \\ &= \underbrace{[(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})]}_{-2\mathbf{H}^\top(\mathbf{y} - \mathbf{H}\mathbf{w})} + \lambda \underbrace{[\mathbf{w}^\top \mathbf{w}]}_{2\mathbf{w}}\end{aligned}$$

Why? By analogy to 1d case...

$\mathbf{w}^\top \mathbf{w}$ analogous to w^2 and derivative of $w^2 = 2w$

Step 3, Approach 1:
Set the gradient = 0

Aside: Refresher on identity matrices

$$I_1 = [1], I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \dots, I_n = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix}$$

Fun facts:

$$\mathbf{I}\mathbf{v} = \mathbf{v}$$

$\stackrel{n \times n}{\mathbf{I}}$ $\stackrel{n \times 1}{\mathbf{v}}$
 \uparrow vector

$$\mathbf{I}\mathbf{A} = \mathbf{A}$$

$\stackrel{n \times n}{\mathbf{I}}$ $\stackrel{n \times n}{\mathbf{A}}$
 \uparrow matrix
 $\mathbf{A}\mathbf{I} = \mathbf{A}$

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

$\stackrel{n \times n}{\mathbf{A}^{-1}}$ $\stackrel{n \times n}{\mathbf{A}}$
 \uparrow by definition
of matrix inverse

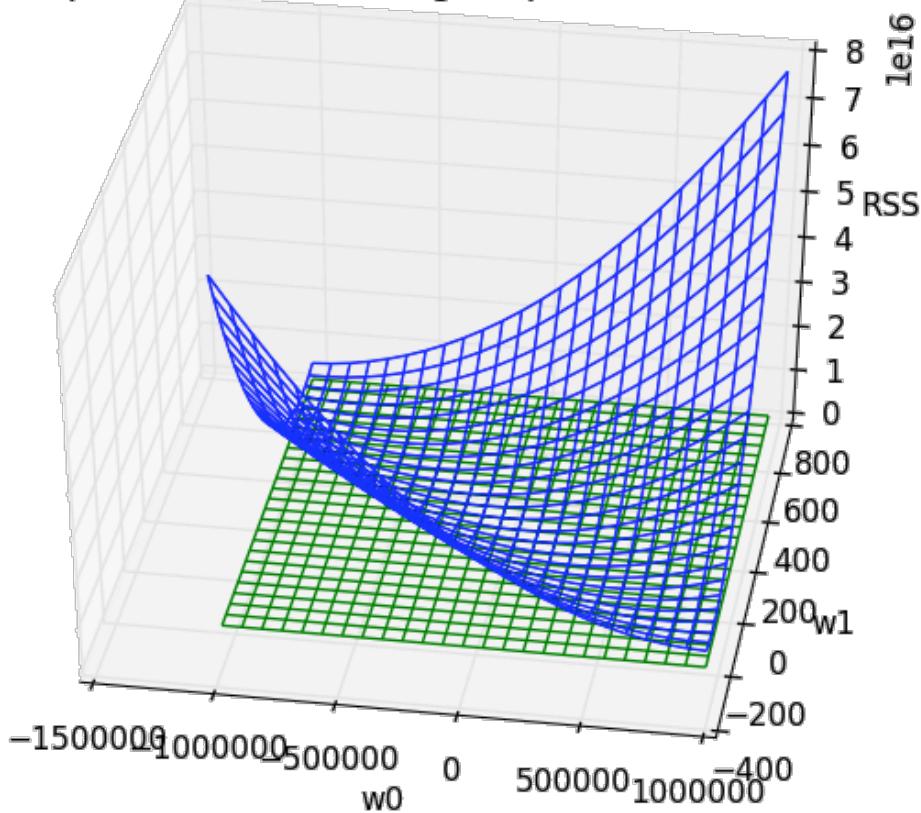
$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$$

$\stackrel{n \times n}{\mathbf{A}}$ $\stackrel{n \times n}{\mathbf{A}^{-1}}$
 \uparrow $\mathbf{A} = \mathbf{A}^{-1}$

$$\begin{aligned} \nabla \text{cost}(\mathbf{w}) &= -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda \boxed{\mathbf{w}} \quad \text{equivalent} \\ &= -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda \boxed{\mathbf{I}\mathbf{w}} \end{aligned}$$

Ridge closed-form solution

3D plot of RSS with tangent plane at minimum



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

Solve for \mathbf{w} :

$$\mathbf{H}^T \mathbf{y} + \mathbf{H}^T \mathbf{H} \hat{\mathbf{w}} + \lambda \mathbf{I} \hat{\mathbf{w}} = 0$$

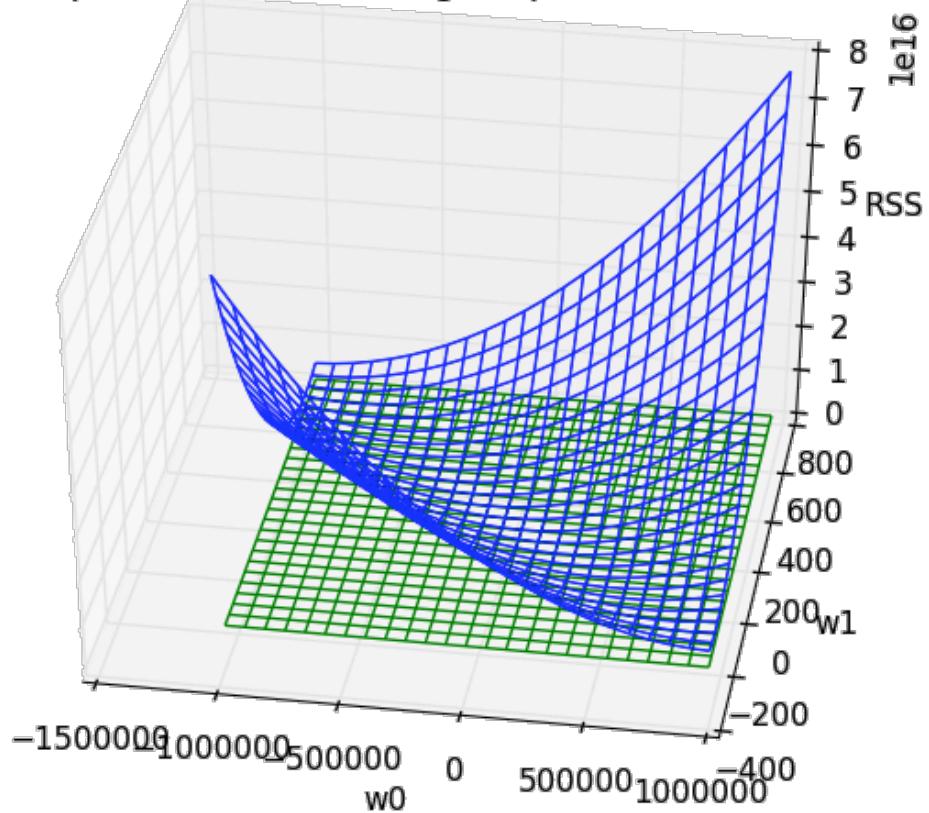
$$\mathbf{H}^T \mathbf{H} \hat{\mathbf{w}} + \lambda \mathbf{I} \hat{\mathbf{w}} = \mathbf{H}^T \mathbf{y}$$

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

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

Interpreting ridge closed-form solution

3D plot of RSS with tangent plane at minimum

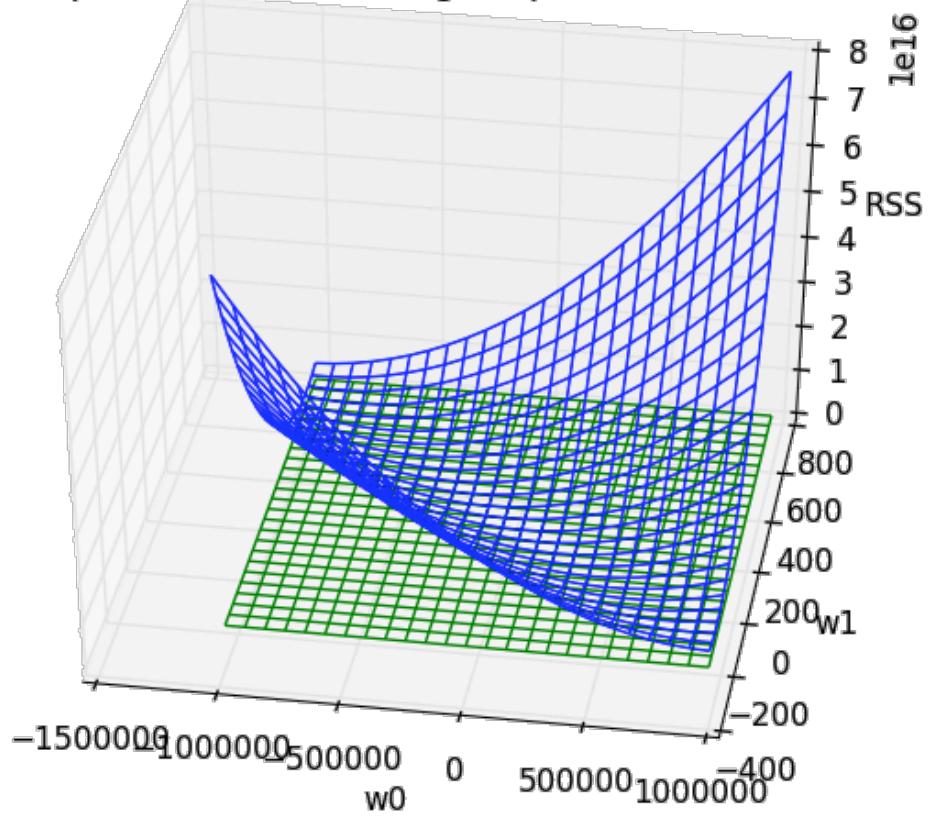


$$\hat{\mathbf{w}}^{\text{ridge}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}$$

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

Interpreting ridge closed-form solution

3D plot of RSS with tangent plane at minimum



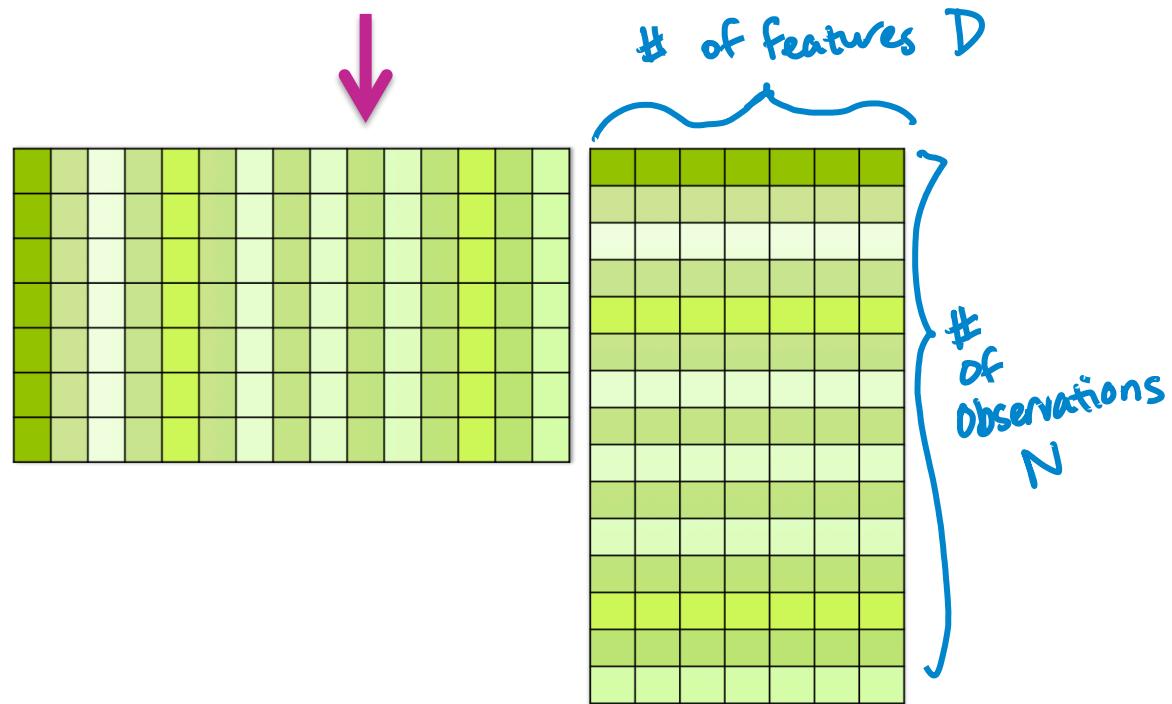
$$\hat{\mathbf{w}}^{\text{ridge}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}$$

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

If $\lambda = \infty$: $\hat{\mathbf{w}}^{\text{ridge}} = 0$ ← because it's like dividing by ∞

Recall discussion on previous closed-form solution

$$\hat{\mathbf{w}}^{\text{LS}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$



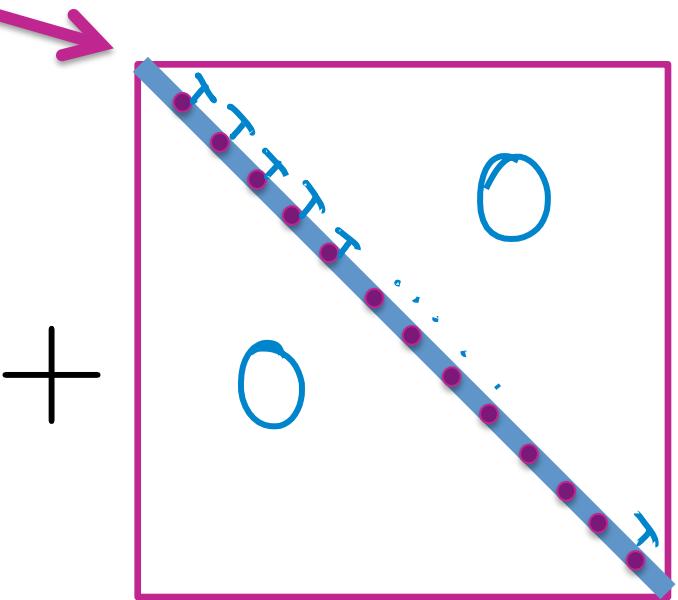
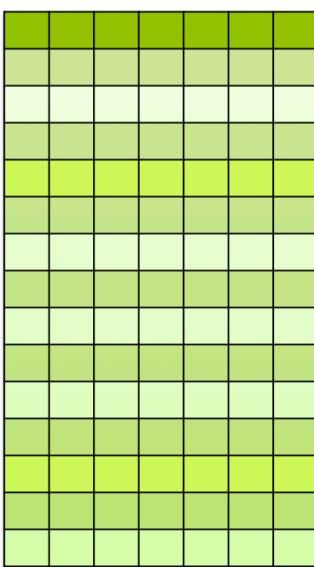
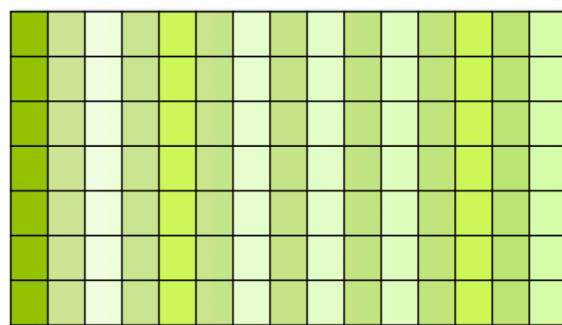
Invertible if:

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

Complexity of inverse:
 $O(D^3)$

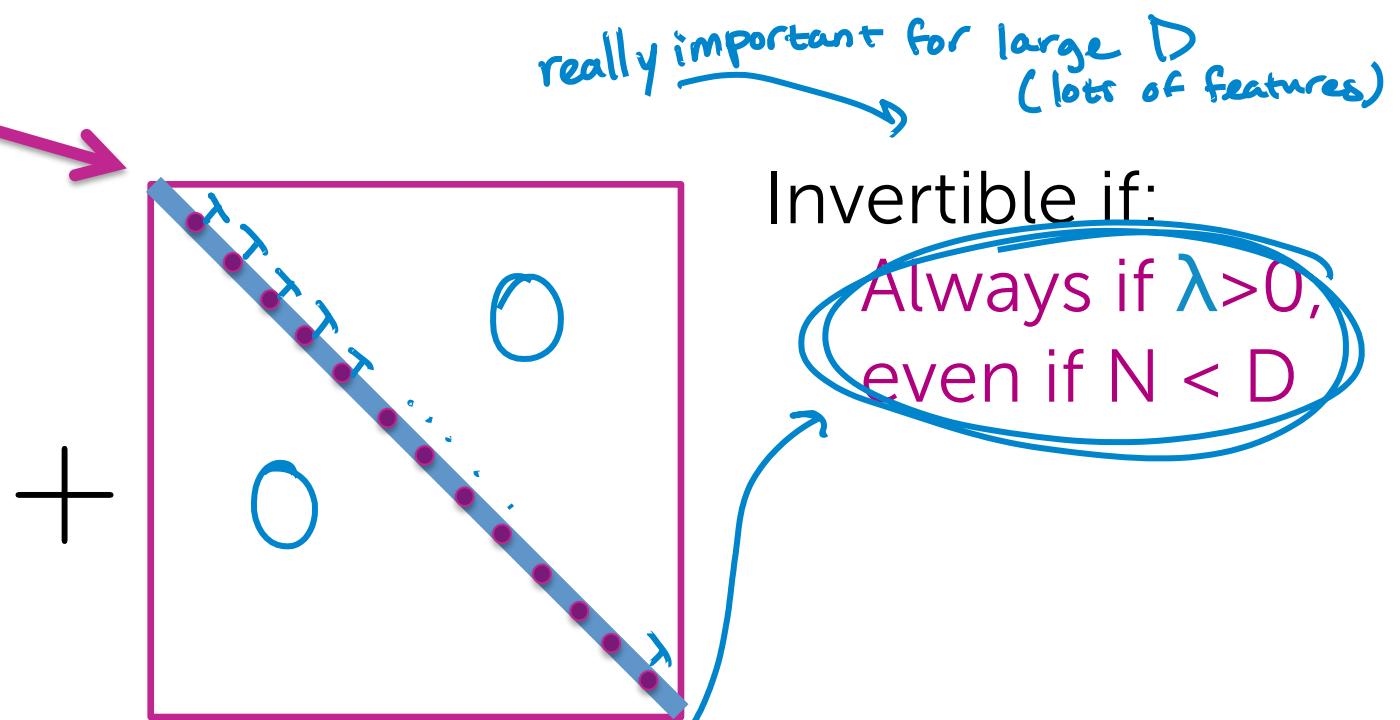
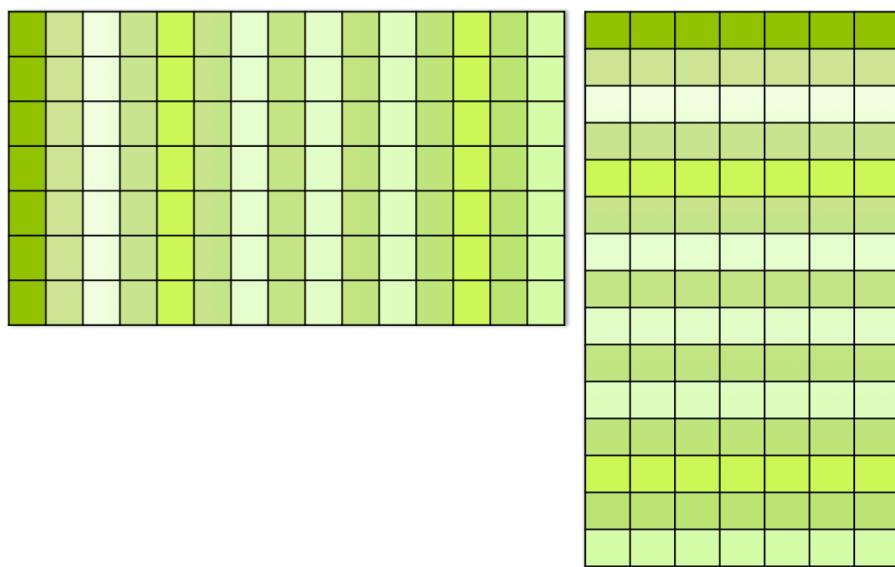
Discussion of ridge closed-form solution

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



Discussion of ridge closed-form solution

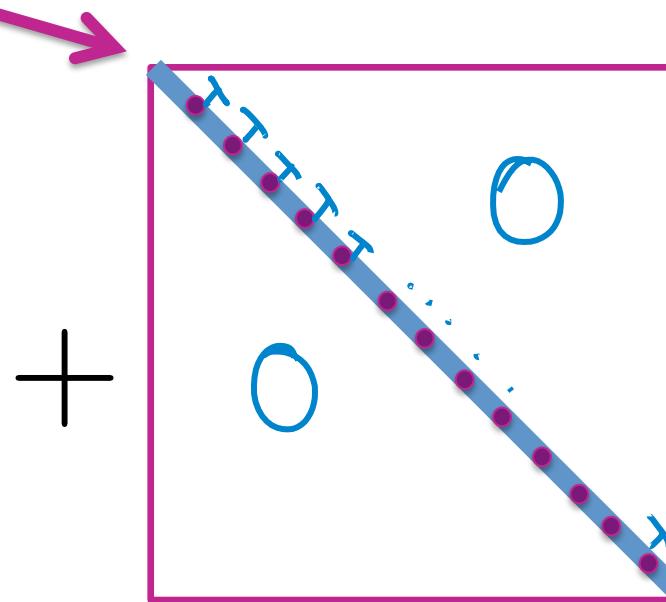
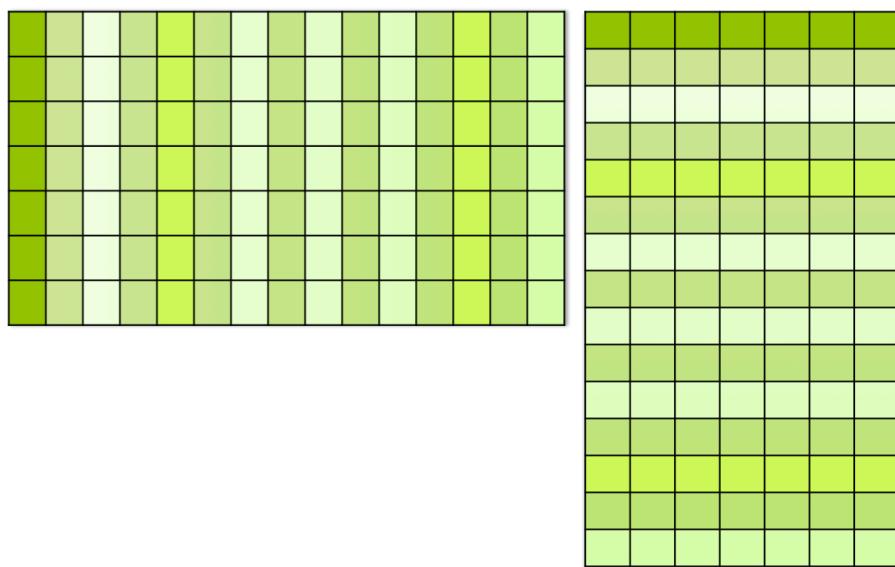
$$\hat{\mathbf{w}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}$$



$\lambda \mathbf{I}$ is making $\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}$ more "regular"
→ "regularized"

Discussion of ridge closed-form solution

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



really important for large D
(lots of features)

Invertible if:

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

Complexity of
inverse:

$O(D^3)$...

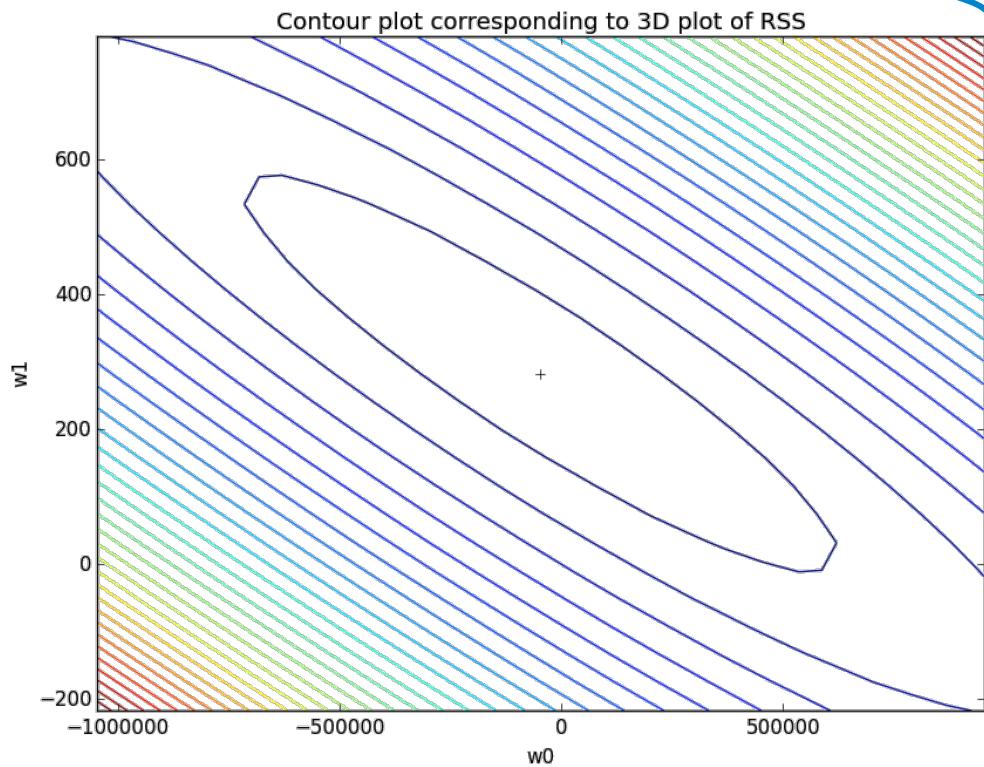
big for large D !

$\lambda \mathbf{I}$ is making
 $\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}$ more "regular"
"regularized"

Step 3, Approach 2: Gradient descent

Elementwise ridge regression gradient descent algorithm

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



Update to j^{th} feature weight:

$$w_j^{(t+1)} \leftarrow \underline{w_j^{(t)}} - \eta *$$

*Same as before
(from RSS term)*

$$\begin{aligned} & \rightarrow 2 \sum_{i=1}^N h_i(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)})) \\ & + 2\lambda w_j^{(t)} \end{aligned}$$

*new term,
comes from the j^{th} component
of $2\lambda\mathbf{w}$*

Elementwise ridge regression gradient descent algorithm

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

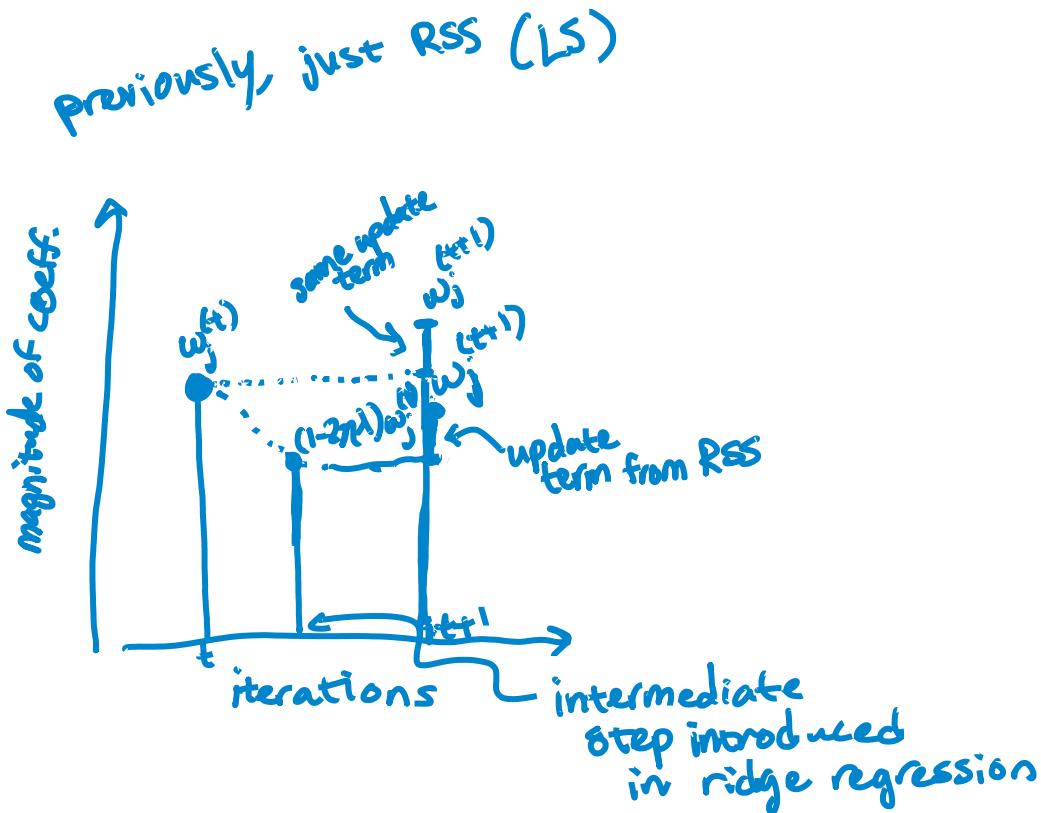
Equivalently:

$$\mathbf{w}_j^{(t+1)} \leftarrow \underbrace{\left(1 - \frac{2n\lambda}{N}\right)}_{\leq 1} \mathbf{w}_j^{(t)} + 2n \sum_{i=1}^N \mathbf{x}_i (\mathbf{y}_i - \hat{\mathbf{y}}_i(\mathbf{w}^{(t)}))$$

increment term

Elementwise ridge regression gradient descent algorithm

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



Equivalently:

$$w_j^{(t+1)} \leftarrow \underbrace{(1-2n\lambda)}_{\leq 1} w_j^{(t)} + \underbrace{2n \sum_{i=1}^N \mathbf{x}_i (y_i - \hat{y}_i(\mathbf{w}^{(t)}))}_{\text{update term from RSS}}$$

$2n\lambda < 1$

Recall previous algorithm

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

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

for $j=0, \dots, D$

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

$$\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} - \eta \text{partial}[j]$$

$$t \leftarrow t + 1$$

Summary of ridge regression algorithm

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

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

for $j=0, \dots, D$

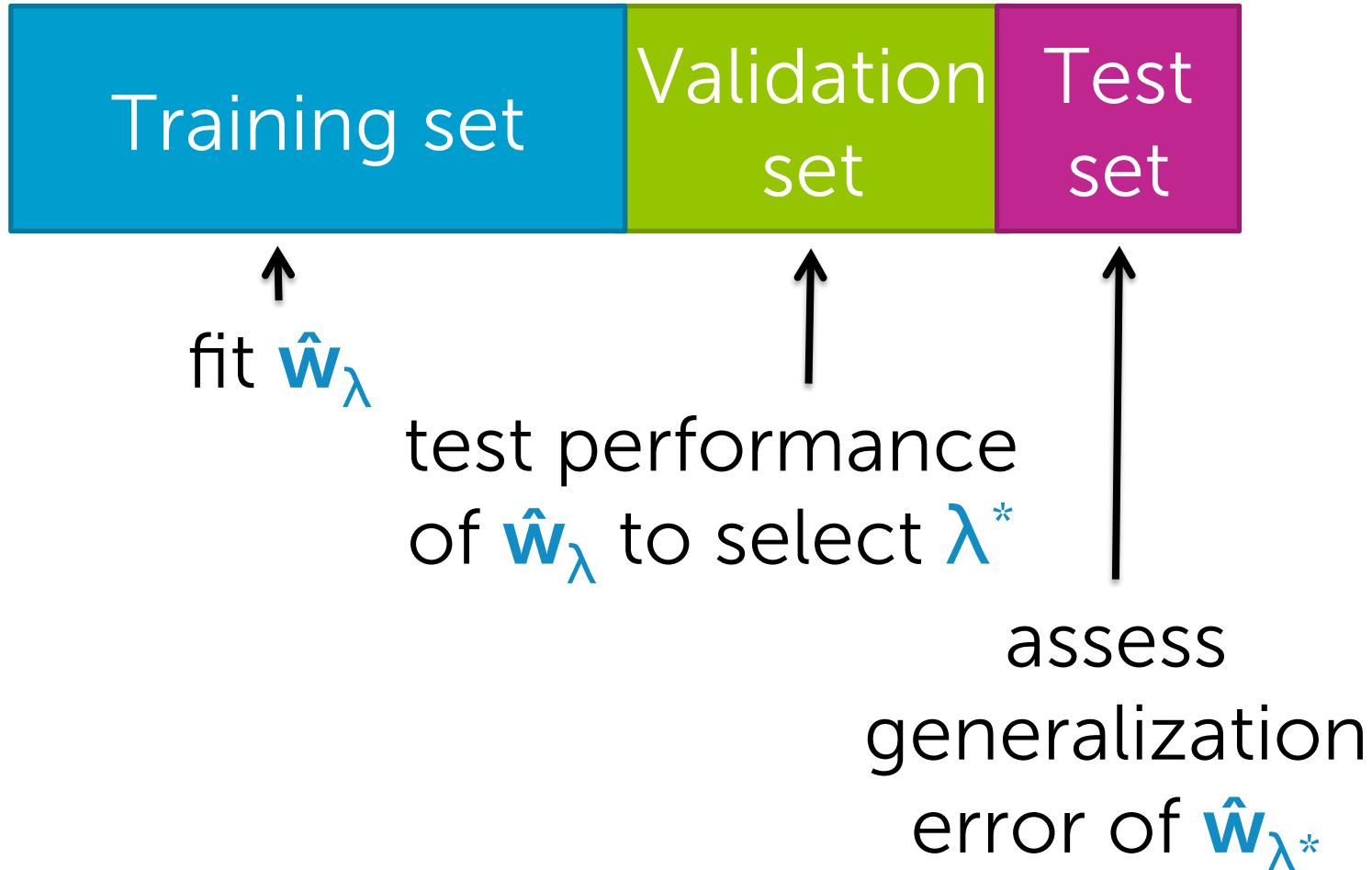
partial[j] = -2 $\sum_{i=1}^N (\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$

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

$t \leftarrow t + 1$

How to choose λ

If sufficient amount of data...



Start with smallish dataset

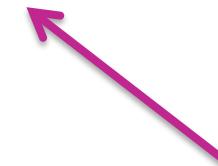
All data

Still form test set and hold out



How do we use the other data?

Rest of data



use for both training and
validation, but not so naively

Recall naïve approach



Is validation set enough to compare performance of \hat{w}_λ across λ values?

No

Choosing the validation set



small validation set

Didn't have to use the last data points
tabulated to form validation set

Can use **any** data subset

Choosing the validation set



Which subset should I use?

Choosing the validation set



Which subset should I use?

ALL!

average
performance
over all
choices

K-fold cross validation

Rest of data

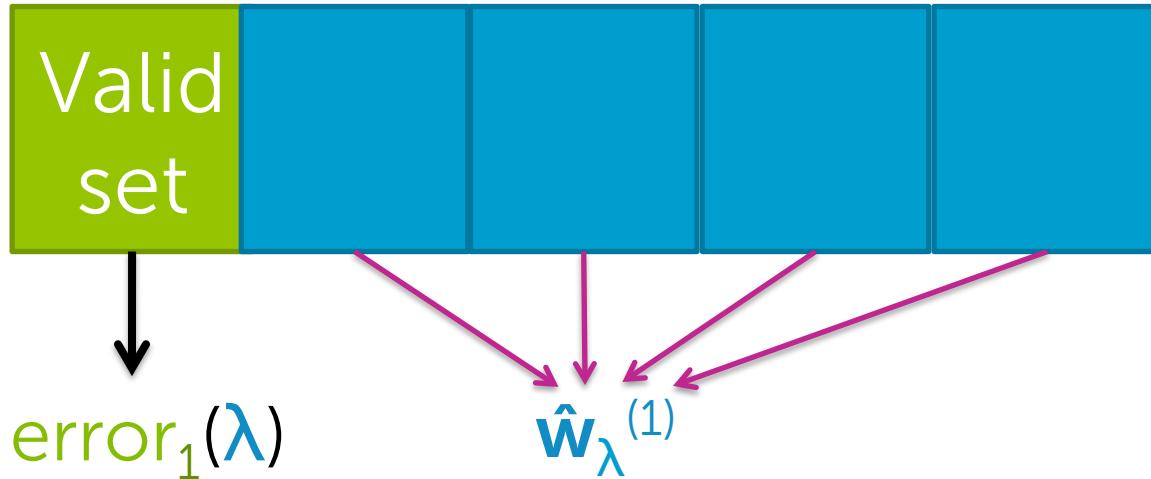
$$\frac{N}{K} \quad \frac{N}{K} \quad \frac{N}{K} \quad \frac{N}{K} \quad \frac{N}{K}$$

Preprocessing:

Randomly assign data to K groups

(use same split of data for all other steps)

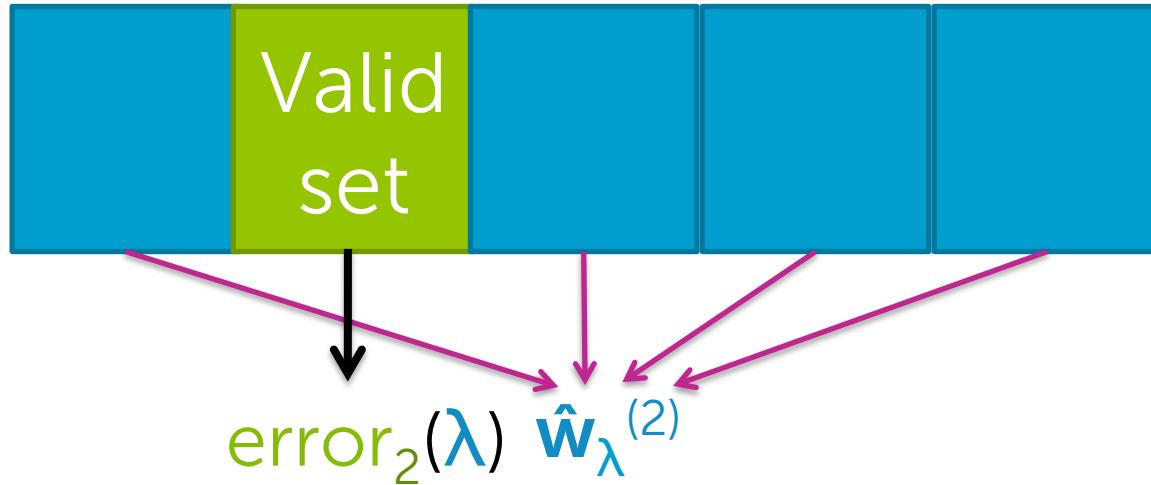
K-fold cross validation



For $k=1, \dots, K$

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

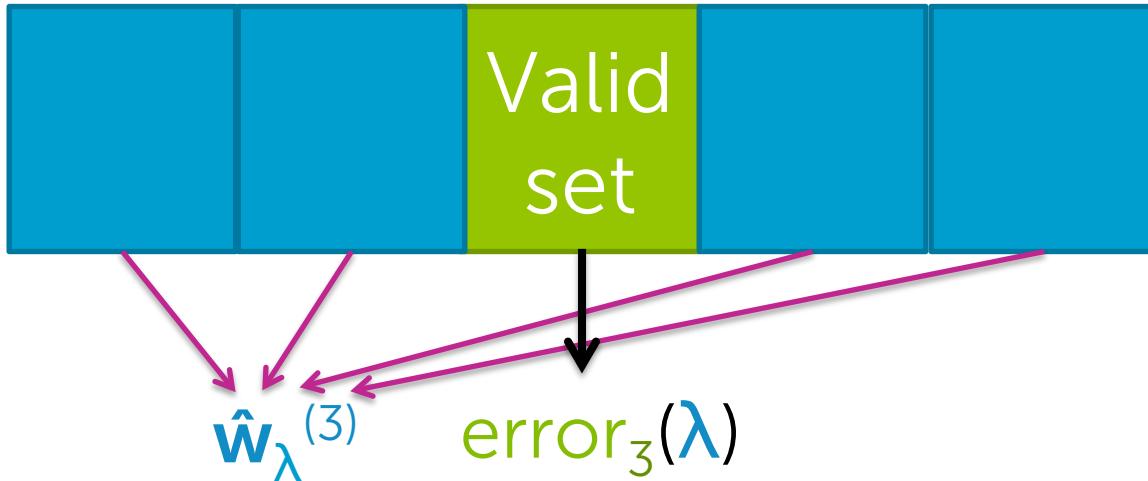
K-fold cross validation



For $k=1, \dots, K$

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

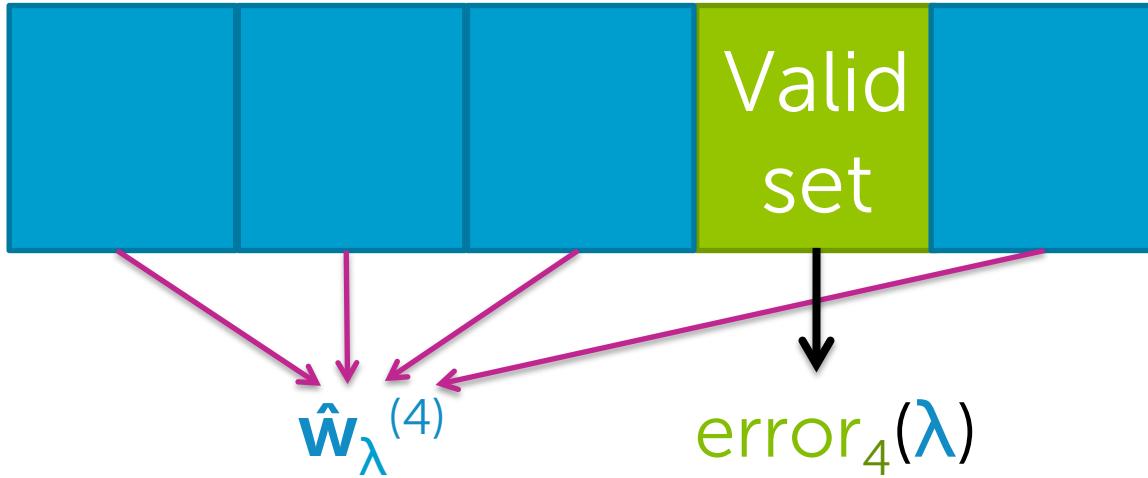
K-fold cross validation



For $k=1, \dots, K$

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

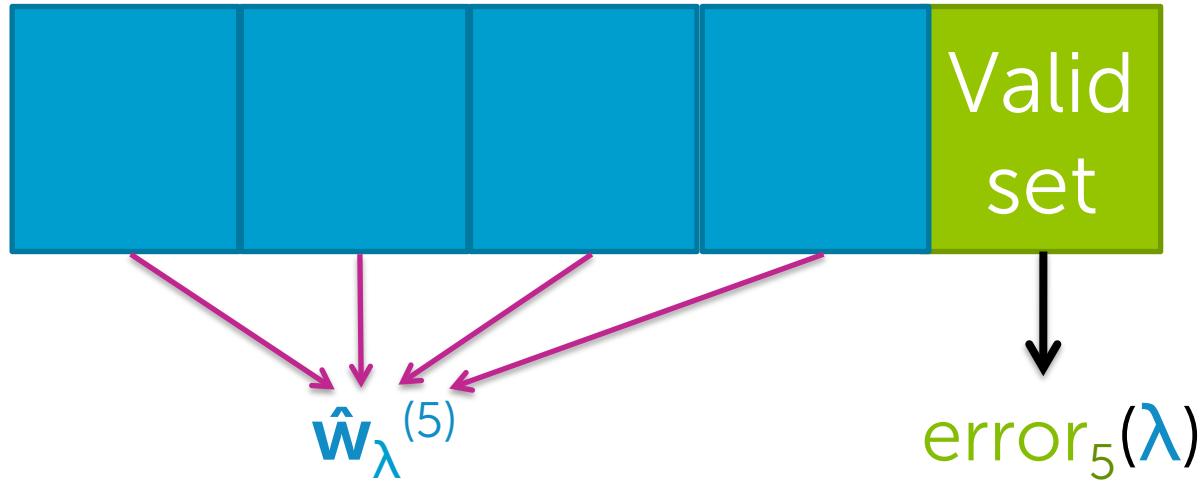
K-fold cross validation



For $k=1, \dots, K$

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

K-fold cross validation

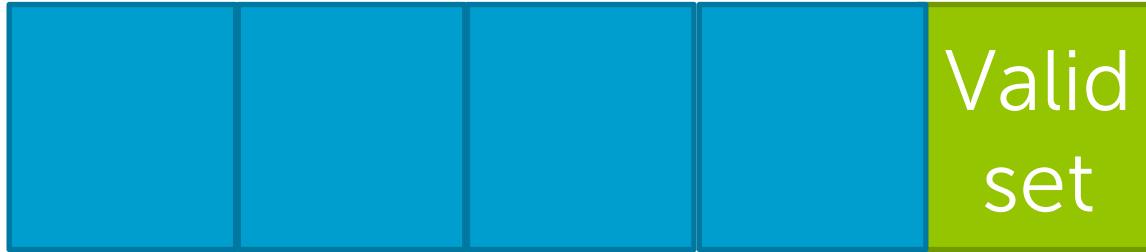


For $k=1, \dots, K$

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

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

K-fold cross validation



Repeat procedure for each choice of λ

Choose λ^* 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 λ

Typically, $K=5$ or 10

5-fold CV

10-fold CV

Summary for ridge regression

What you can do now...

- Describe what happens to magnitude of estimated coefficients when model is overfit
- Motivate form of ridge regression cost function
- Describe what happens to estimated coefficients of ridge regression as tuning parameter λ is varied
- Interpret coefficient path plot
- Estimate ridge regression parameters:
 - In closed form
 - Using an iterative gradient descent algorithm
- Implement K-fold cross validation to select the ridge regression tuning parameter λ