

Ridge Regression

- High complexity model -> **Low Bias** and **High Variance** and vice-versa.
- Bias-Variance trade-off - required to achieve good predictive performance.
- **Ridge Regression** - automatically balance between Bias and Variance.

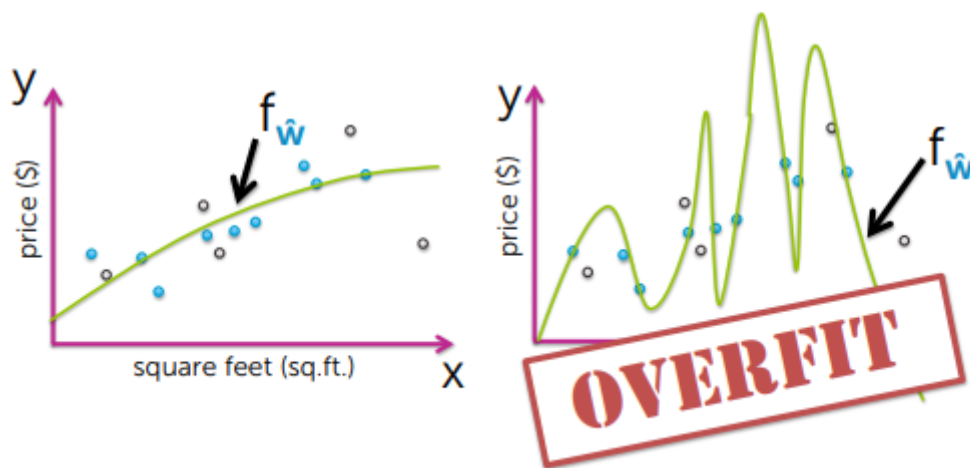
Overfit:

- When a model is highly specific to the training data and does not generalize well - then it is **Overfit**.
- Training error of model w' > training error of another model w_1
- True error of model w' < True error of another model w_1

Polynomial -> features are power of an input.

Flexibility of high-order polynomials

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



When models become overfit the estimated parameters (\hat{w}) become very large in magnitude.

Overfitting is generic issue with complex models

- It occurs with polynomial regression.
- Occurs in complex models
 - lots of inputs (d large).
 - lots of features (D large). $y(j) = \sum_{j=0}^{D-1} w(j)h(j) + \epsilon(i)$

How does # of observations influence overfitting?

- Few observations (**N small**) -> **rapidly overfit** as model **complexity increases**. With few points, as the order of the polynomial increases, it becomes easy to hit all the datapoints - hence overfit the dataset.
- Many observations (**N is large**) -> **harder to overfit**.
- It doesn't overfit easily since there are dense observations to overfit the input. It's not able to hit all the datapoints in the observations.

How does # of inputs influence overfitting?

- 1 input (e.g., sqft):
- The dataset must be very dense - must include representative examples of all possible (sqft, 'dollars') pairs to avoid overfitting.
- This is hard, to have all possible (sqft, 'dollars') pairs values.
- d inputs (e.g., sqft, #bath, #bed, lot size, year,...):
- The data must include examples of all possible (sqft, #bath, #bed, lot size, year,...,'\$') combos to avoid overfitting.
- This is even harder to cover all possible scenarios.

Balancing the fit and magnitude of coefficients

- Overfitting increases the magnitude of the coefficient.

Quality Metrics -> Adding term to cost-of-fit to prefer small coefficients.

- Thus far -> quality metrics - depended on the **actual sales price / actual output** and the **predicted sales price / predicted output** - **RSS** - was used to measure of fit.
- Now The Quality Metrics - is also gonna incorporate the **complexity of the model**.

Desired total cost format

Want to balance :

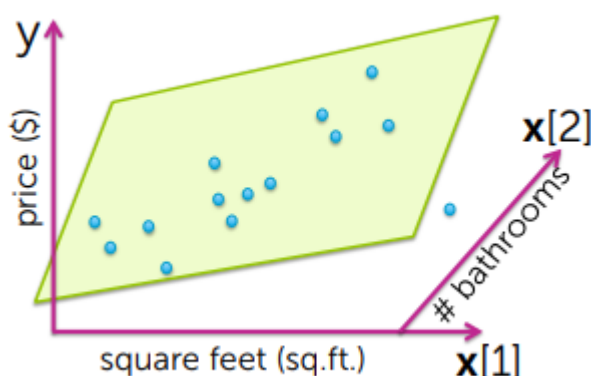
1. How well the function fits the data.
2. The complexity of the model - through 'Magnitude of the estimated coefficients'.

Previous cost = $RSS = (Actual - Predicted)^2$ -> measure of fit

New Cost = Total cost = measure of fit + measure of magnitude of coefficients

- Measure of fit -> (small # = good fit to training data).
- Measure of magnitude of coefficients -> (small # = not overfit).

I. Measure of fit



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

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

pred. value using \mathbf{w}

small RSS \rightarrow model fitting training data well

II. Measure of magnitude of the coefficients

What summary # is indicative of size of regression coefficients?

- Sum? $w_0 = 1,527,301$ $w_1 = -1,605,253$
 $w_0 + w_1 = \text{small } \neq$
- Sum of absolute value?
 $|w_0| + |w_1| + \dots + |w_D| = \sum_{j=0}^D |w_j| \triangleq \|w\|_1$ L_1 norm ... discuss more in next module
- Sum of squares (L_2 norm)
 $w_0^2 + w_1^2 + \dots + w_D^2 = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2$ L_2 norm ... focus of this module

Total Cost

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

Resulting Objective

- Find the estimated parameters (\hat{w}) that minimizes the **total cost**.
- A tuning parameter λ is introduced in order to balance of fit and magnitude.

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

λ tuning parameter = balance of fit and magnitude

If $\lambda = 0$:
 reduces to minimizing $\text{RSS}(w)$, as before (old solution) $\rightarrow \hat{w}^{\text{LS}}$ ← 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 $0 \leq \|\hat{w}\|_2 \leq \|\hat{w}^{\text{LS}}\|_2$

Ridge Regression - a.k.a $L(2)$ regularization = $\text{RSS}(w) + \lambda \|w\|_2^2$

Bias-variance tradeoff

Large λ : low complex model

high bias, low variance

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

In essence, λ
controls model
complexity

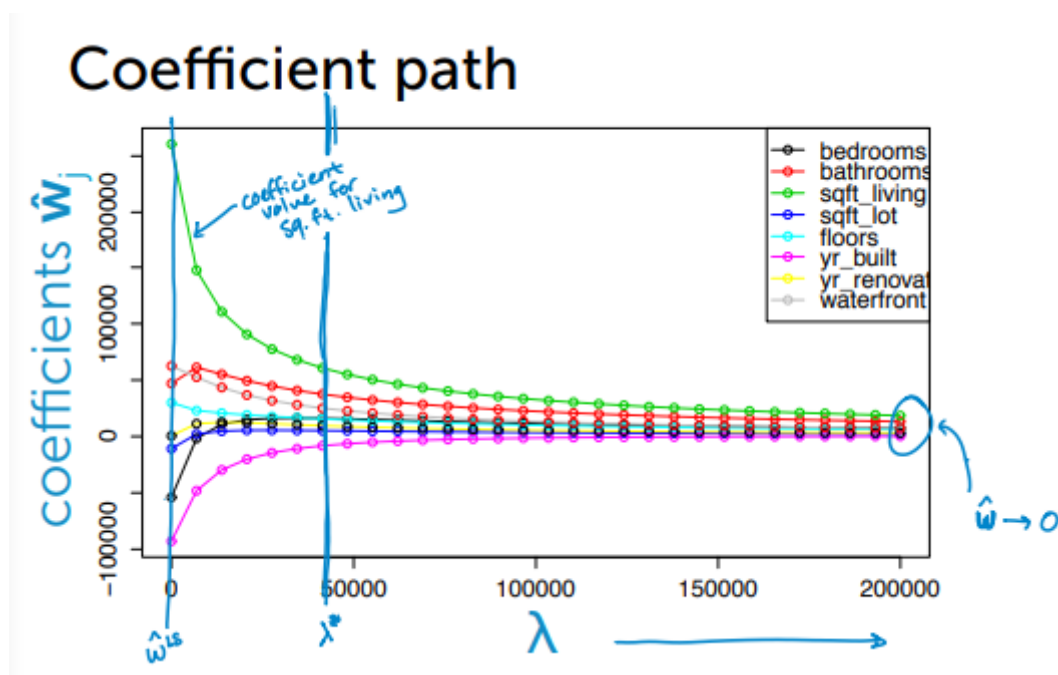
Small λ : high complex model

low bias, high variance

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

The ridge coefficient path:

- In general when λ is small, the coefficient magnitude is large.
- When λ is large \rightarrow infinity, the coefficient magnitude is small $\rightarrow 0$;



ML Algorithm \rightarrow ML Block Diagram

Computing the gradient of the Ridge Objective

Step 1: Rewrite total cost in matrix notation

- Model for all N observations together. $\mathbf{Y} = \mathbf{H} * \mathbf{w} + \epsilon$

Recall matrix form of RSS

Model for all N observations together

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

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

Rewrite magnitude of coefficients in vector notation

$$\begin{aligned} \|\mathbf{w}\|_2^2 &= w_0^2 + w_1^2 + w_2^2 + \dots + w_D^2 \\ &= \begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_D \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \\ &= \mathbf{w}^T \mathbf{w} \end{aligned}$$

Putting it all together

In matrix form, ridge regression cost is:

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

Step 2: Compute the gradient

Gradient of ridge regression cost

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

Why? By analogy to 1d case...

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

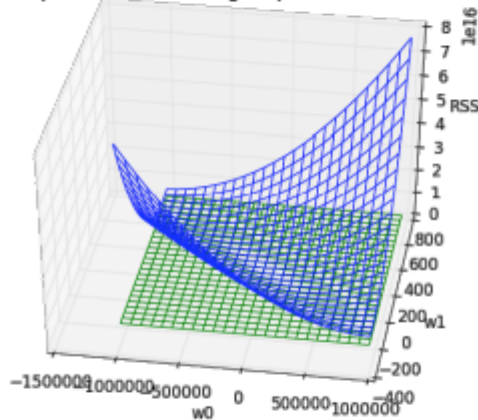
Step 3: Approach 1 : Set the gradient = 0 | Closed-form solution

Ridge closed-form solution

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

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

3D plot of RSS with tangent plane at minimum



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

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

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

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

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

Interpreting ridge closed-form solution

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

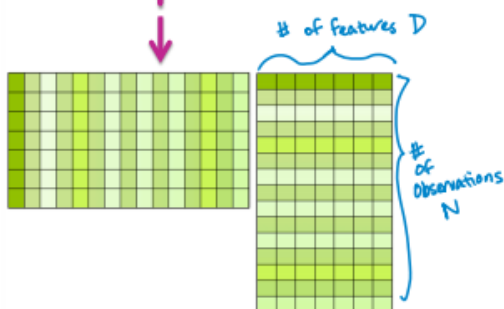
$$\text{If } \lambda = 0: \hat{\mathbf{w}}^{\text{ridge}} = (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{y} = \hat{\mathbf{w}}^{\text{LS}} \leftarrow \text{old solution!}$$

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

Closed-form solution : w.r.t Ridge regression -> the λ term is multiplied with \mathbf{I} (identity matrix) -> λ is a scalar -> identity matrix with λ along the diagonal and the rest are 0 is formed.

previous closed-form solution | ridge closed-form solution

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



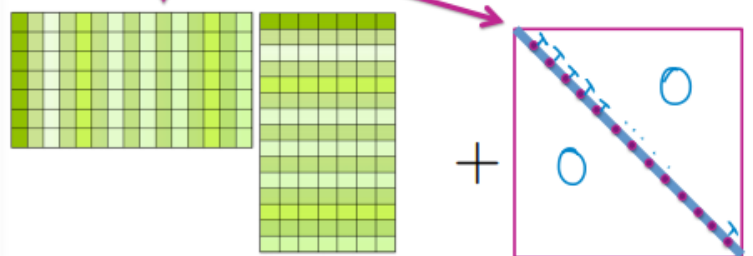
Invertible if:

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

Complexity of inverse:

$O(D^3)$

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



Invertible if:

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

really important for large D
(lots of features)

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"

- The closed-form solution is computationally expensive.

Step 3: Approach 2 : Gradient Descent

- The new coefficient

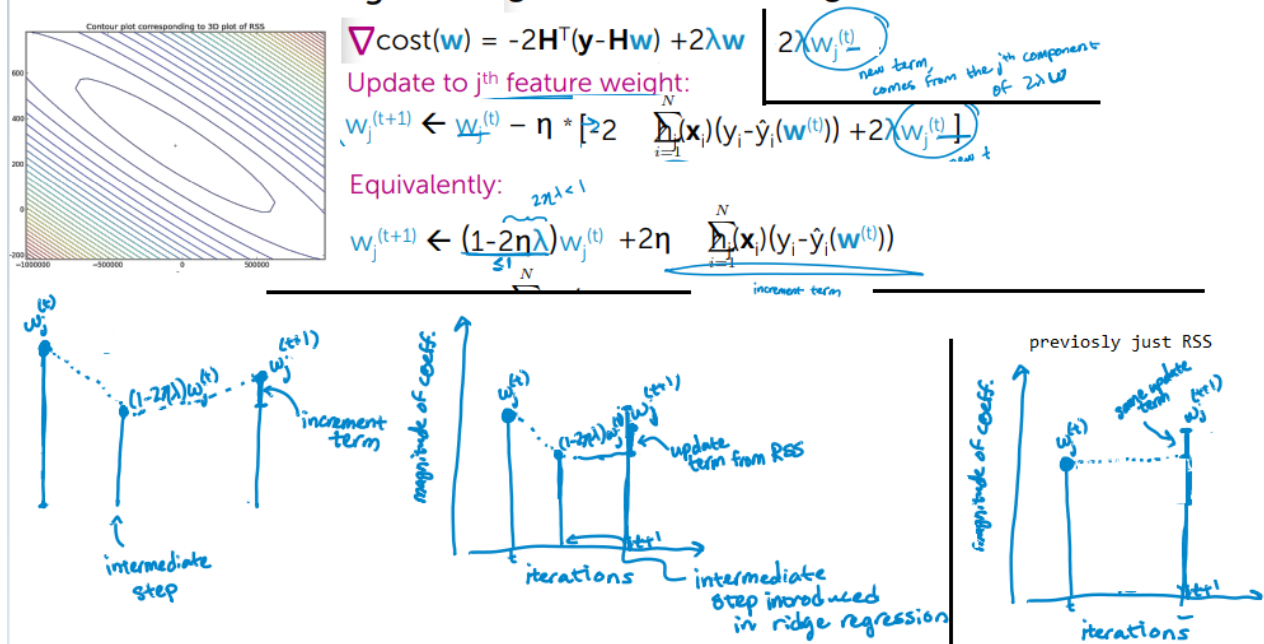
$$w(j)(t+1) = (1-2\eta\lambda) * w(j)(t) + 2\eta \sum_{i=1}^N h_j(x) (y_i - \hat{y}_i(w))$$

- Step i: With every new iteration $\rightarrow (1-2\eta\lambda)$ reduces the $w(j)(t)$ term since η & λ are > 0 .
- Step ii: At $t+1$ iteration, the update term for RSS is added to the previously shrink term. This is the new $w(j)(t+1)$ - coefficient.

Previously - Just RSS

- Here the $w(j)(t)$ term was taken and for the new $w(j)(t+1)$ term the $w(j)(t)$ was added with **update term** from RSS.

Elementwise ridge regression gradient descent algorithm



Gradient Descent - implemented w.r.t 'Multiple Regression' v/s 'Ridge Regression'

previous algorithm

```

init  $w^{(1)} = 0$  (or randomly, or smartly),  $t = 1$ 
while  $\|\nabla \text{RSS}(w^{(t)})\| > \epsilon$ 
  for  $j = 0, \dots, D$ 
    partial[j] =  $-2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$ 
     $w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \text{partial[j]}$ 
   $t \leftarrow t + 1$ 

```

ridge regression algorithm

```

init  $w^{(1)} = 0$  (or randomly, or smartly),  $t = 1$ 
while  $\|\nabla \text{RSS}(w^{(t)})\| > \epsilon$ 
  for  $j = 0, \dots, D$ 
    partial[j] =  $-2 \sum_{i=1}^N h_j(x_i)(y_i - \hat{y}_i(w^{(t)}))$ 
     $w_j^{(t+1)} \leftarrow (1-2\eta\lambda)w_j^{(t)} - \eta \text{partial[j]}$ 
   $t \leftarrow t + 1$ 

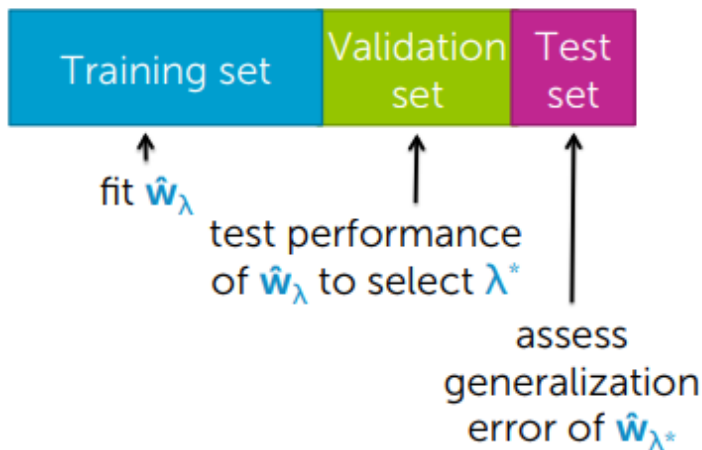
```

Selecting tuning parameters λ via cross validation

Case 1 : Sufficient amount of data

Practical implementation

1. Select λ^* such that $\hat{\mathbf{w}}_{\lambda^*}$ minimizes error on validation set
2. Approximate generalization error of $\hat{\mathbf{w}}_{\lambda^*}$ using test set



Typical splits

Training set	Validation set	Test set
80%	10%	10%
50%	25%	25%

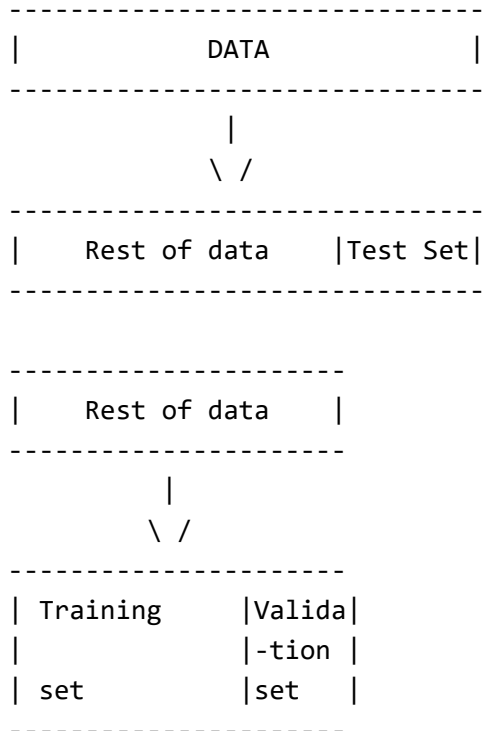
$\hat{\mathbf{w}}_{\lambda}$ = estimate parameters on training data

λ^* = tuning parameter to control the model complexity with lowest test error

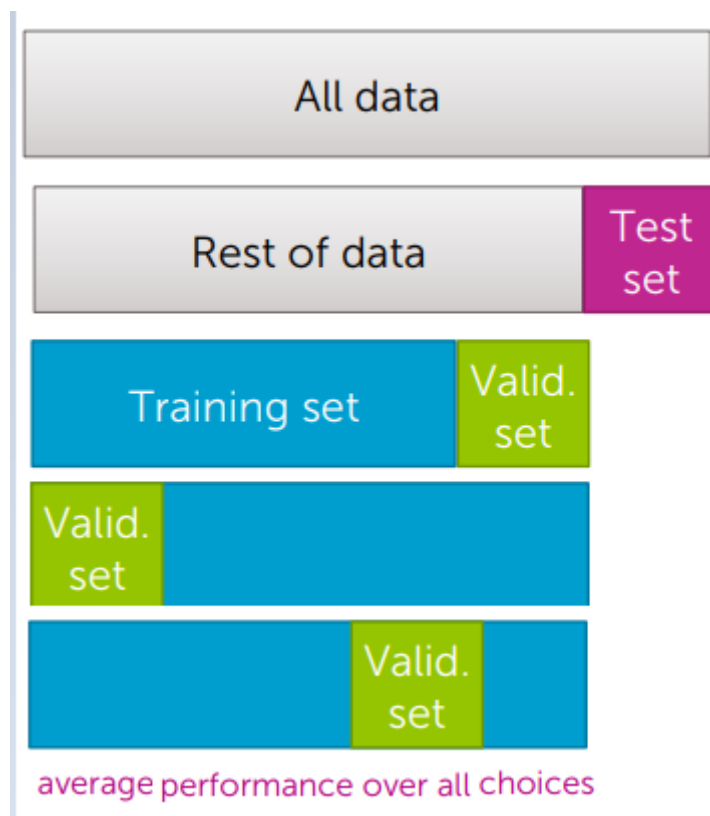
$\hat{\mathbf{w}}_{\lambda^*}$ = Fitted model for selected complexity λ^*

- For every value of that tuning parameter we can fit our model in the training data.
- Then assess the performance of the fitted model on a validation set, and can tabulate all values of 'lambda' that we can consider and choose the specific model complexity according to error on validation set.
- Assess the performance of the selected model on the test set.
- In the presence of **sufficient data** the above splitting process can be applied even to 'Ridge Regression'.

Case 1 : Insufficient amount of data



- Is the validation set enough to compare performance of estimated parameters for a given tuning parameter 'lambda' ($\hat{\mathbf{w}}$) across all 'lambda' λ values?
- It is not necessary to use the last data points tabulated to form the validation set. Rather can use any subset.
- In case of insufficient data. Can use the entire training dataset in subsets. **Then Average the performance over all choices of the subset.**



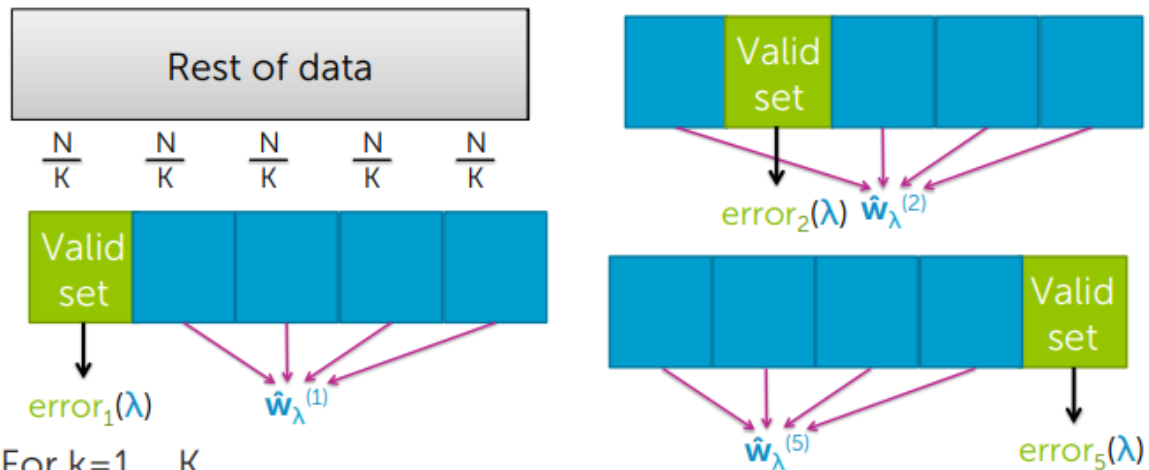
K-fold cross validation

- Step 1 : Preprocessing : Randomly assign the data to K groups.

- Take the Rest of the data (data apart from the test set) and divide into into k groups. There are N total observations so every block is gonna have N/K observations and are assigned randomly to each block.
- Step 2: For $k = 1, \dots, K$
 - i. Estimate $(\hat{w})_{\lambda(k)}$ on the training blocks (except the k -subset);
 - ii. Compute the error on the validation block (K) : $\text{error}_k(\lambda)$
- Step 3: Compute the average error : $\text{CV}(\lambda) = 1/K \sum_{k=1}^K \text{error}_k(\lambda)$.
- Step 4: Repeat the above procedure for each choice of λ .
 - Choose λ^* to minimize $\text{CV}(\lambda)$.

K-fold cross validation

Preprocessing: Randomly assign data to K groups



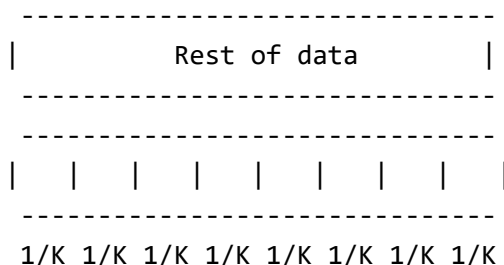
1. Estimate $\hat{w}_{\lambda}^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

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

Choose λ^* to minimize $\text{CV}(\lambda)$

What value of K?

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



- This is the **leave-out-one cross validation** (LOO Cross Validation).
- It is computationally intensive - requires computing N fits of model per λ .
- Typically , $K = 5$ (**5-fold CV**) or $k = 10$ (**10-fold CV**).

How to handle the intercept

- w_0 -> first coefficient of the model is usually the intercept, i.e where x is 0. $h_0(x) = \text{constant}$ (1).
- Therefore the H matrix first column will be 1.
- Since w_0 is the intercept.
- The cost of Ridge Regression -> $\text{RSS}(w) + \lambda \|w\|^2$;
 - where λ - strength of penalty.
 - since w_0 is the intercept multiplied by 1, squaring it will encourage the intercept w_0 to also be small.
 - This operation will not be indicative of overfitting. (Overfitting is indicated by the magnitude of the coefficients.)

Avoid the intercept term from being squared.

Option 1: Don't penalize the intercept

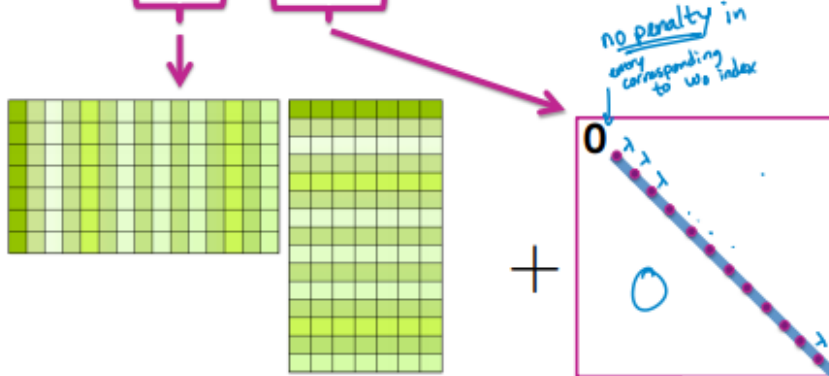
- Previous Ridge Regression Cost : $\text{RSS}(w) + \lambda \|w\|^2$;
- New Ridge Regression Cost : $\text{RSS}(w_0, w_{\text{rest}}) + \lambda \|w_{\text{rest}}\|^2$;

Closed-form solution

- Here the Identity matrix ' I ' gets modified to ' I^{mod} ' - the I^{mod} matrix has the first element as 0, and the rest of the diagonal is filled with λ and the rest of the matrix is filled with 0's.

Option 1: Don't penalize intercept – Closed-form solution –

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



new penalty:
 $\lambda w_{\text{rest}}^T w_{\text{rest}}$

gradient:

$$2\lambda w_{\text{rest}} = 2\lambda I w_{\text{rest}} \rightarrow 2\lambda \begin{bmatrix} 0 & & \\ & 1 & \\ & & \ddots \\ & & & 1 \end{bmatrix} \begin{bmatrix} w_0 \\ w_{\text{rest}} \end{bmatrix}$$

Gradient Descent algorithm

- While performing the update, in case $j = 0$, w_0 (intercept - coefficient) - no shrinking operation is performed, only the fit to data - **update term from RSS is added.*
- For all other cases -> the previous coefficient is shrunk and then appended to form the next coefficient of the polynomial regression term.

Option 1: Don't penalize intercept

– Gradient descent algorithm –

```

while || $\nabla$  RSS( $\mathbf{w}^{(t)}$ )|| >  $\epsilon$ 
  for j=0,...,D
    partial[j] = -2  $\sum_{i=1}^N (\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$ 
    if j==0
       $w_0^{(t+1)} \leftarrow w_0^{(t)} - \eta \text{ partial[j]}$   $\leftarrow$  old LS update (no shrinkage to  $w_0$ )
    else  $\leftarrow$  for all other features
       $w_j^{(t+1)} \leftarrow (1-2\eta\lambda)w_j^{(t)} - \eta \text{ partial[j]}$   $\leftarrow$  ridge update
  t  $\leftarrow$  t + 1

```

Option 2 : Center the 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 descent algorithm).

Quiz

1
point

1. Which of the following is NOT a valid measure of overfitting?

- ☒ Sum of parameters ($w_1 + w_2 + \dots + w_n$)
- ☐ Sum of squares of parameters ($w_1^2 + w_2^2 + \dots + w_n^2$)
- ☐ Range of parameters, i.e., difference between maximum and minimum parameters
- ☐ Sum of absolute values of parameters ($|w_1| + |w_2| + \dots + |w_n|$)

- Explanation : Take weight vector in 2D like the following one : $\mathbf{w} =$ From the value of the weights, we might suspect overfitting. Now: If you consider the sum of weights, you get -200. So this is not a good measure for overfitting detection. On the other hand, if you consider range, you get $2500 - (-2700) = 5200$, which seems good measure for overfitting detection.

1
point

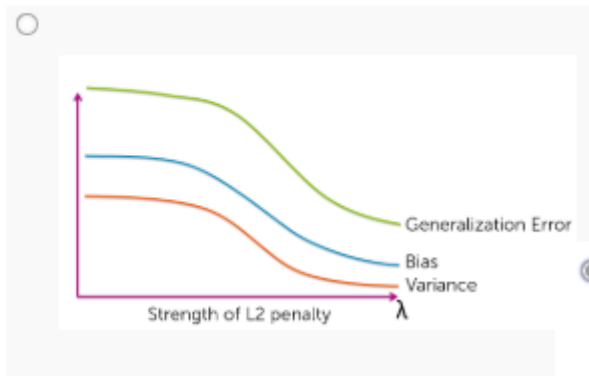
2. In ridge regression, choosing a large penalty strength λ tends to lead to a model with (choose all that apply):

- ☒ High bias
- ☐ Low bias
- ☐ High variance
- ☒ Low variance

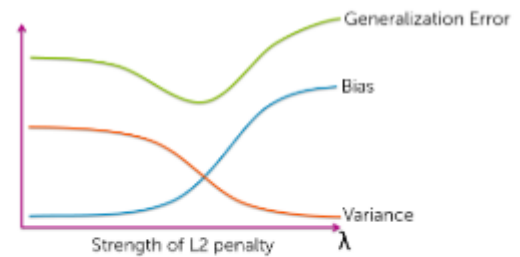
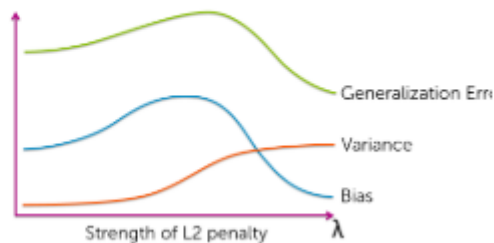
*** When lambda is large (model complexity is low) - Bias is high and variance is low.**

1 point

3. Which of the following plots best characterize the trend of bias, variance, and generalization error (all plotted over λ)?



☐



- The bias increases with 'lambda'.
- The variance decreases with 'lambda'.
- The Generalization error does not proportionally vary with model complexity. It decreases when the bias meets the variance. Increases elsewhere.

1 point

4. In ridge regression using unnormalized features, if you double the value of a given feature (i.e., a specific column of the feature matrix), what happens to the estimated coefficients for every other feature? They:

- ☐ Double
- ☐ Half
- ☐ Stay the same
- ☒ Impossible to tell from the information provided

- Suppose, $w_0 + w_1 x_1 + w_2 x_2 = y$.
- Let $x_1 = 2$ and $x_2 = 4$ for $y = 7$, then one set of likely values of (w_0, w_1, w_2) will be $(1, 1, 1)$.
- If $x_2 = 8$, w_2 will have to be half of previous value and will be less than other two parameters.
- Since ridge regression penalizes larger parameters more, using lambda, values of other parameters should change.
- Hence the estimated parameters may or may not change.

1
point

5. If we only have a small number of observations, K-fold cross validation provides a better estimate of the generalization error than the validation set method.

- ☒ True
- ☐ False

1
point

6. 10-fold cross validation is more computationally intensive than leave-one-out (LOO) cross validation.

- ☐ True
- ☒ False

1
point

7. Assume you have a training dataset consisting of N observations and D features. You use the closed-form solution to fit a multiple linear regression model using ridge regression. To choose the penalty strength λ , you run leave-one-out (LOO) cross validation searching over L values of λ . Let $\text{Cost}(N, D)$ be the computational cost of running ridge regression with N data points and D features. Assume the prediction cost is negligible compared to the computational cost of training the model. Which of the following represents the computational cost of your LOO cross validation procedure?

- ☐ $LN \cdot \text{Cost}(N, D)$
- ☒ $LN \cdot \text{Cost}(N - 1, D)$
- ☐ $LD \cdot \text{Cost}(N - 1, D)$
- ☐ $LD \cdot \text{Cost}(N, D)$
- ☐ $L \cdot \text{Cost}(N - 1, D)$
- ☐ $L \cdot \text{Cost}(N, D)$

- Since it is LOO, each block is 1 observation wide. therefore 1 observation is a validation set , while the rest (N-1) observations are a part of the training set.
- Therefore computational cost = $L N \text{Cost}(N-1, D)$

1
point

8. Assume you have a training dataset consisting of 1 million observations. Suppose running the closed-form solution to fit a multiple linear regression model using ridge regression on this data takes 1 second. Suppose you want to choose the penalty strength λ by searching over 100 possible values. How long will it take to run leave-one-out (LOO) cross-validation for this selection task?

- ☐ About 3 hours
- ☐ About 3 days
- ☒ About 3 years
- ☐ About 3 decades

- $N = \# \text{ of observation} = 10^6$;
- $L = \# \text{ of } \lambda = 100$;
- $\text{Cost}(N-1, D) \text{ operation} = 1 \text{ second}$
- $T = L N \text{Cost}(N-1, D) = 100 \cdot 10^6 \cdot 1 \text{ second}$
- $T = 100 \cdot 10^6 \cdot 1 \text{ second} / (60 \cdot 60 \cdot 24 \cdot 365) \sim 3 \text{ years.}$

1
point

9. Assume you have a training dataset consisting of 1 million observations. Suppose running the closed-form solution to fit a multiple linear regression model using ridge regression on this data takes 1 second. Suppose you want to choose the penalty strength λ by searching over 100 possible values. If you only want to spend about 1 hour to select λ , what value of k should you use for k -fold cross-validation?

- ☐ $k=6$
- ☒ $k=36$
- ☐ $k=600$
- ☐ $k=3600$

- k - folds cross validaions.
- $L = \# \text{ lambda} = 100$;
- $t = \text{time to select lambda* over all lambda} = 1 \text{ hour}$
- $k = t / L = 1 \text{ hour} / 100 = 60 * 60 / 100 = 36$;

In []: