

Math 7243

Machine Learning and Statistical Learning Theory

### Section 3. Ridge and LASSO Regressions

Instructor: He Wang

Department of Mathematics  
Northeastern University

## Review :

Training Data  $D = \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(n)}, y^{(n)})\}$

1. Linear regression
2. Residual sum of squares
3. Matrix calculus

Model Assumption:  $h(\vec{x}) = \vec{\theta}^T \vec{x} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$

Find  $\vec{\theta}$  to minimize  $RSS(h) = \sum_{i=1}^n (h(\vec{x}^{(i)}) - y^{(i)})^2$

*if full column  
rank*

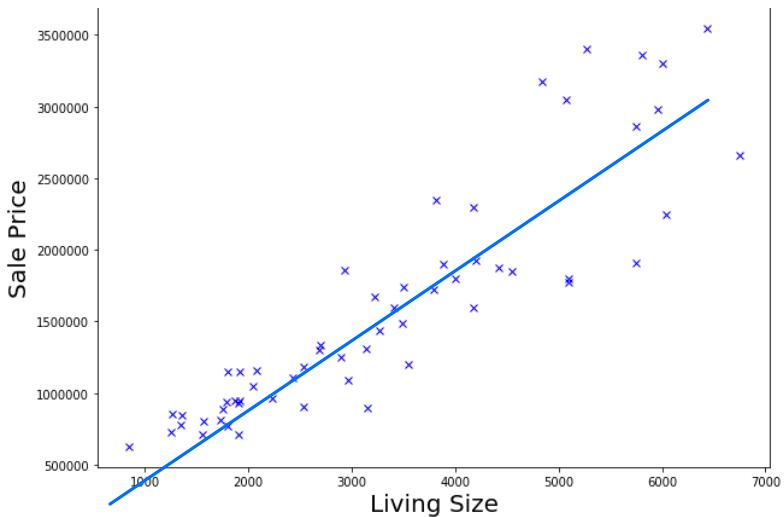
Least Squares solution  $\vec{\theta} = (X^T X)^{-1} X^T \vec{y}$

## Today:

1. Locally weighted linear regression
2. Interpretation in Probability
3. Ridge Regression
4. Lasso Regression
5. Elastic net Regression

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

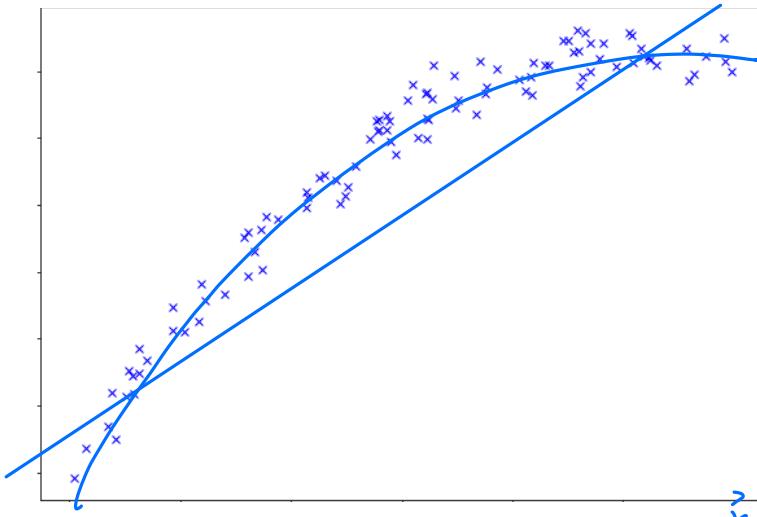
➤ Predict house price



Wrong assumption

**Potential Disadvantage of a parametric approach:** The (linear) model we choose will usually not match the true unknown form of  $h(\vec{x})$ . If the chosen model is too far from the true  $h(\vec{x})$ , then our estimate will be poor.

We can try to solve this problem by choosing flexible models that can fit many different possible functional forms flexible for  $h(\vec{x})$ . But in general, fitting a more flexible model requires estimating a greater number of parameters.

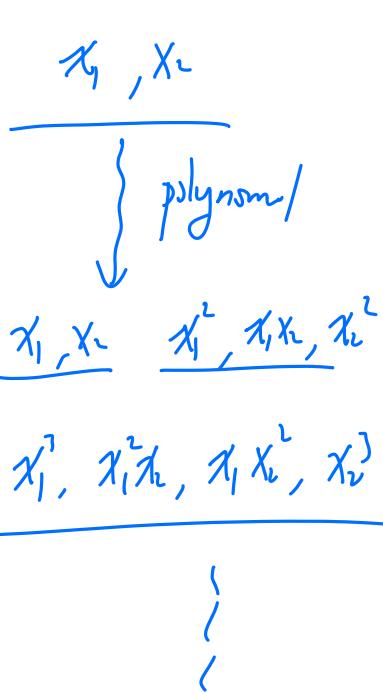


New assumption :

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 \sqrt{x} + \theta_4 \log(x)$$

Data

$x_1$	$y$	$x_2$	$x_1$	$x_4$
$x_1^{(1)}$	$y^{(1)}$	$x^{(1)2}$	$\sqrt{x^{(1)}}$	$\log(x^{(1)})$
$x_1^{(2)}$	$y^{(2)}$	$x^{(2)2}$	$\sqrt{x^{(2)}}$	$\log(x^{(2)})$
.	.	.	.	.
1	1	1	1	1



Potential Disadvantage: overfitting the data.

These more complex models can lead to a phenomenon that they follow the errors, or noise, too closely.

- Locally weighted regression  
(non-parametric model)

Goal: Evaluate  $h$  at certain  $x$

- Recall Linear Regression: Find

$$h(\vec{x}) = \vec{\theta}^T \vec{x} = \theta_0 + \theta_1 x_1$$

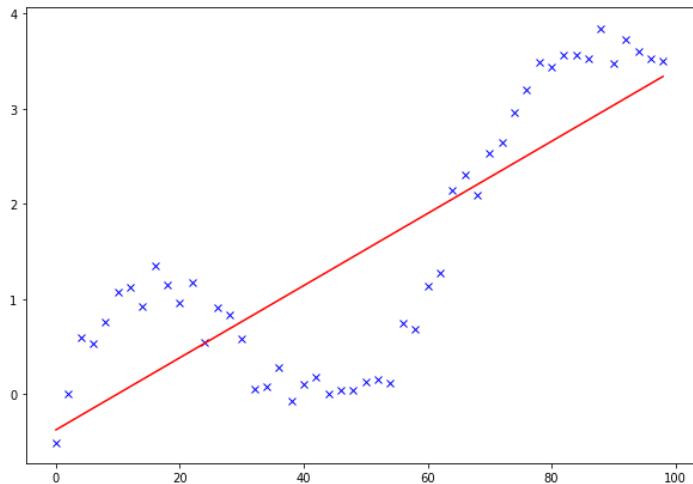
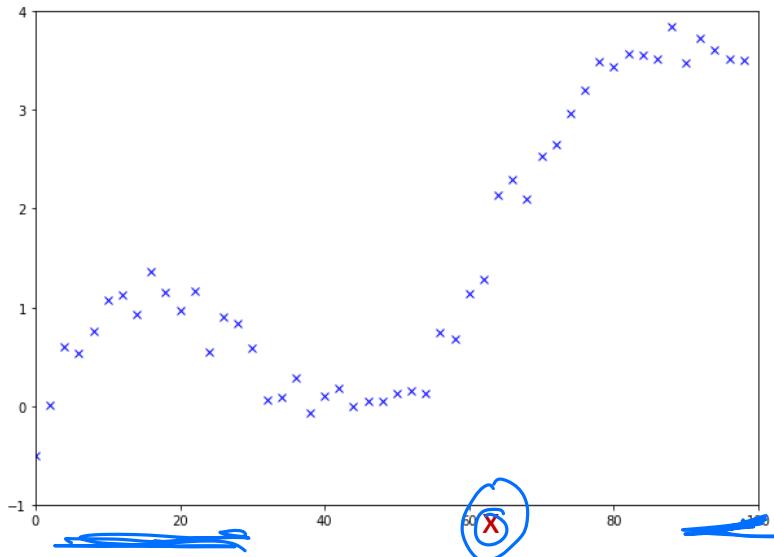
to minimize the cost function:

$$RSS(\vec{\theta}) = \sum_{i=1}^n (h(\vec{x}^{(i)}) - \vec{y}^{(i)})^2$$

Add weight to the data

$$J(\vec{\theta}) = \sum_{i=1}^n w^{(i)} (\vec{\theta}^T \vec{x}^{(i)} - y^{(i)})^2$$

minimize

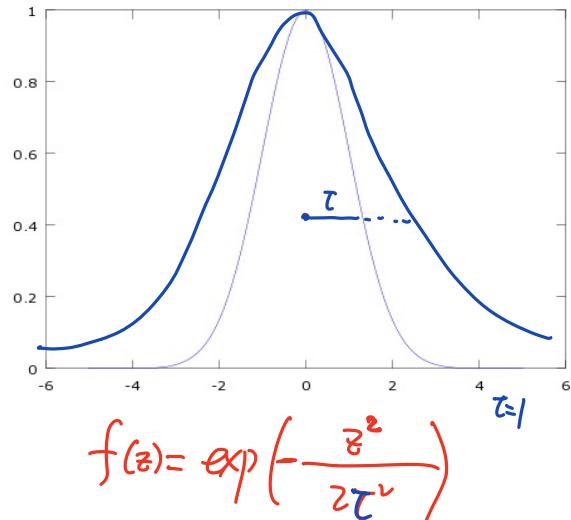


$$w^{(i)} = \exp \left( -\frac{\|\vec{x}^{(i)} - \vec{x}\|^2}{2\tau^2} \right)$$

► Minimize new cost function

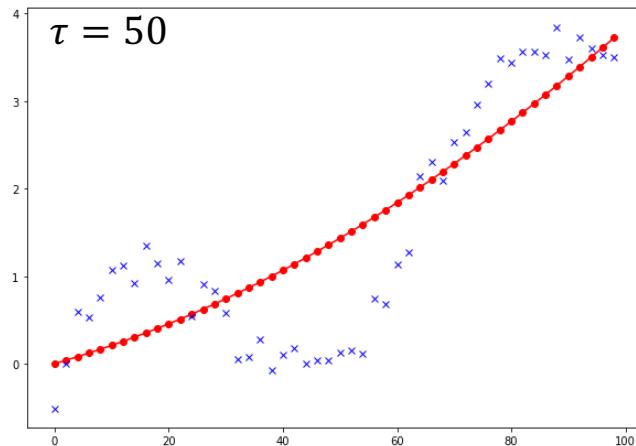
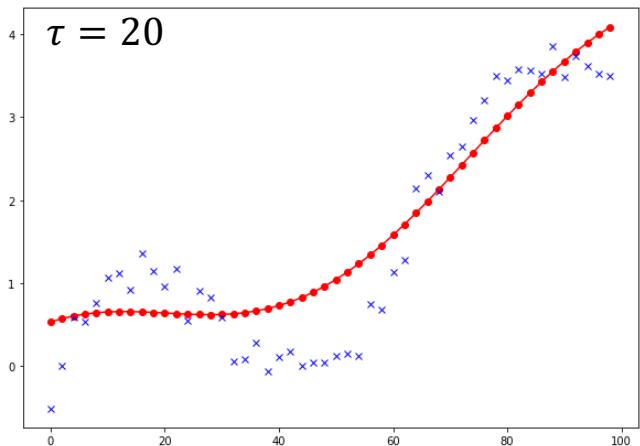
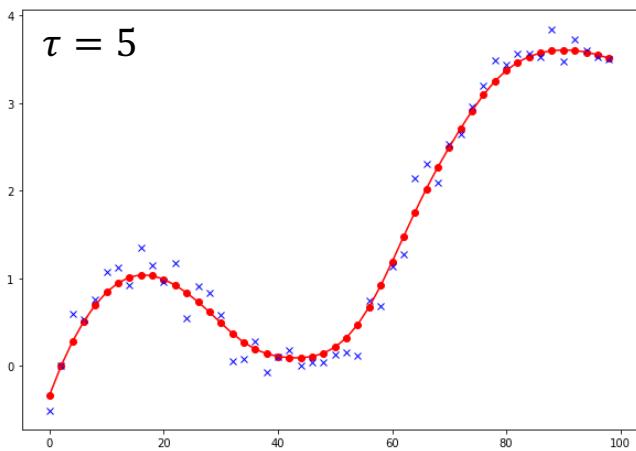
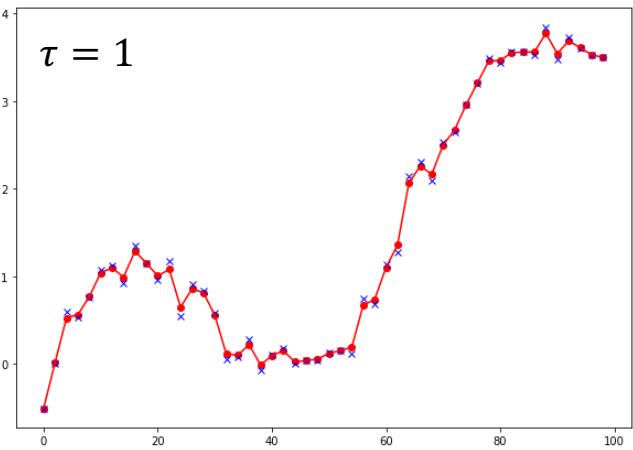
$$w^{(i)} \rightarrow \begin{cases} 1 & \text{if } \|\vec{x}^{(i)} - \vec{x}\| \rightarrow 0 \\ 0 & \text{if } \|\vec{x}^{(i)} - \vec{x}\| \rightarrow \infty \end{cases}$$

$$J(\vec{\theta}) = (\vec{X}\vec{\theta} - \vec{y})^T W (\vec{X}\vec{\theta} - \vec{y})$$



• Claim:  $\nabla_{\vec{\theta}} J(\vec{\theta}) = 2 \vec{X}^T W (\vec{X}\vec{\theta} - \vec{y}) \underset{=0}{=} 0$        $W = \begin{bmatrix} w^{(1)} & & & \\ & w^{(2)} & & 0 \\ 0 & \ddots & \ddots & w^{(n)} \end{bmatrix}$

$$\vec{\theta} = \underbrace{(\vec{X}^T W \vec{X})^{-1}}_{\text{---}} (\vec{X}^T W \vec{y})$$



We need the training data as well as the parameters to make a prediction.

```
import numpy as np → Vectors  
import pandas as pd → data
```

```
# To plot pretty figures  
%matplotlib inline  
import matplotlib as mpl  
import matplotlib.pyplot as plt
```

```
x = 2*np.arange(50)  
y = np.sin(x/10) + (x/50)**2 + 0.2*np.random.randn(50)
```

```
def normal_equation(x, y, w=None):  
    if w is None:  
        return np.linalg.inv(x.T.dot(x)).dot(x.T).dot(y)  
    else:  
        return np.linalg.inv(x.T.dot(w).dot(x)).dot(x.T).dot(w).dot(y)|
```

$$(X^T X)^{-1} X^T \bar{y})$$

$$(X^T w X)^{-1} (X^T w \bar{y})$$

```
def weight_w(x, x_i, tau):  
    return np.diag(np.exp(-((x-x_i)[:, 1]**2)/(2*tau**2)))
```

```
xa=np.append(np.ones(x.shape[0]), x)  
xc = xa.reshape(2,50).T  
yc = y.reshape(50,1)
```

```
theta = normal_equation(xc, yc)|
```

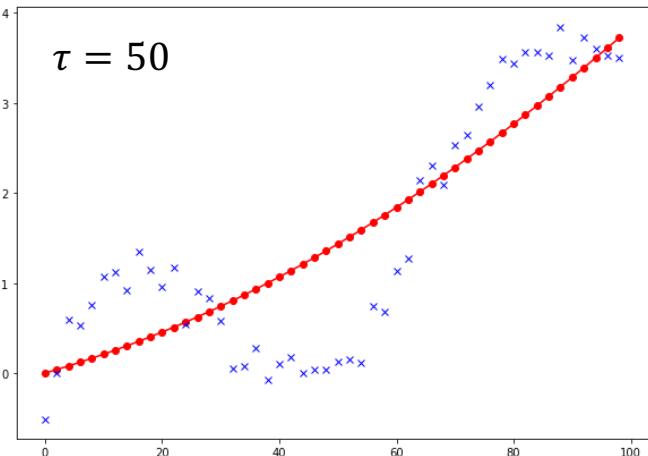
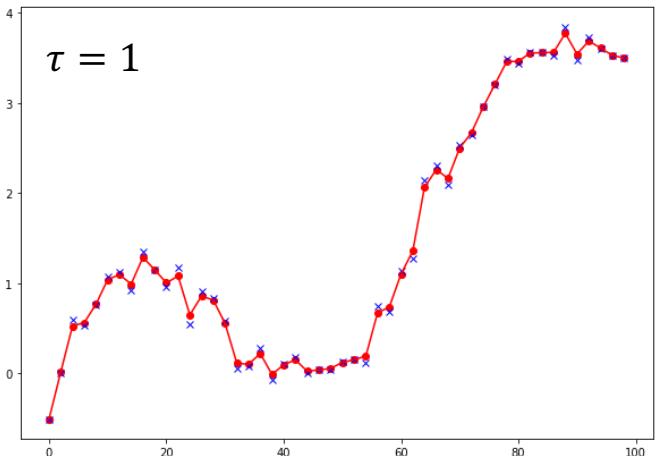
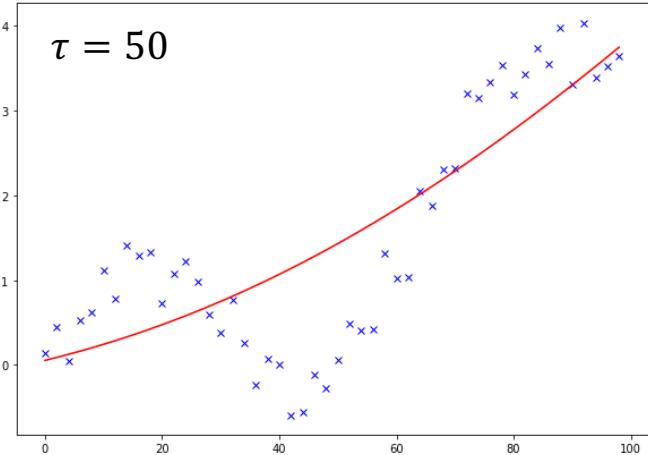
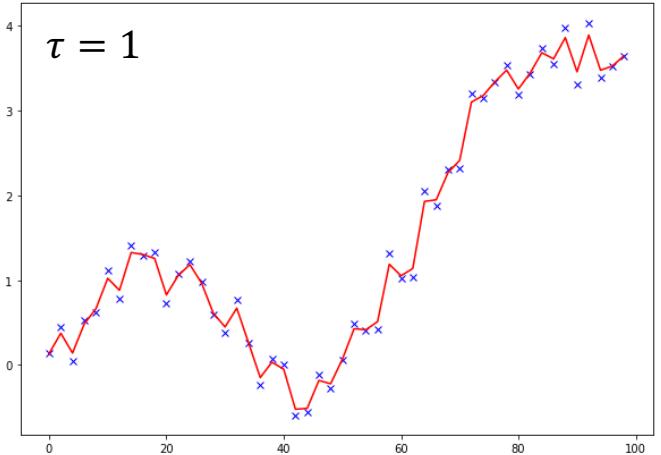
$$X = \begin{bmatrix} & \\ & \\ & \end{bmatrix}$$

```
fig, ax = plt.subplots()
Y_pred=xc.dot(theta)
ax.plot(xc[:,1],yc,'x',color='Blue')
plt.plot(xc[:,1], Y_pred, color='red')
fig.set_size_inches(10, 7)
plt.show()
```

→ least square

```
# Calculate predictions
pred = []
for k, xc_j in enumerate(xc):
    w = weight_w(xc, xc_j, 5)
    theta = normal_equation(xc, yc, w)
    pred.append(theta.T.dot(xc_j[:,np.newaxis]).ravel()[0])
```

```
fig, ax = plt.subplots()
ax.plot(xc[:,1],yc,'x',color='Blue')
plt.plot(xc[:,1], pred, color='red')
plt.scatter(xc[:,1], pred, color ='red')
fig.set_size_inches(10, 7)
plt.show()
```



## ➤ Bias-Variance Trade Off

When  $\tau = 1$ , if we change the training data from the same model and recompute the linear classifier, the change in the fit is very high. *high variance.*

By contrast, when  $\tau = 50$ , if we change the training data from the same model and recompute the linear classifier, the change in the fit is very low. *low variance high bias*

This contrast between the two algorithms is known as the **bias-variance trade off**.

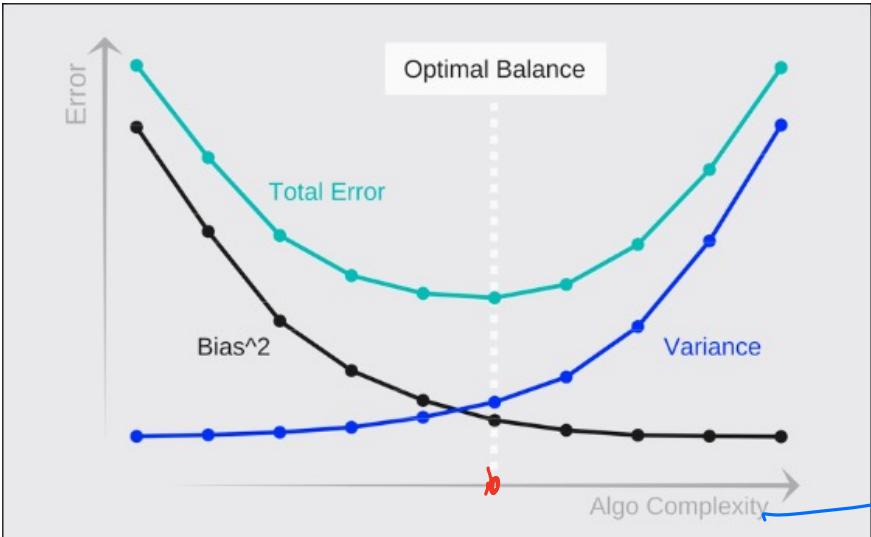
For a class of models, the **bias** roughly is the expected error of the best classifier in the model class given a random set of training data. (This part of the generalization error is due to wrong assumptions.)

The **variance** is roughly the sensitivity of the model to the training data. (This part is due to the model's excessive sensitivity to small variations in the training data.)

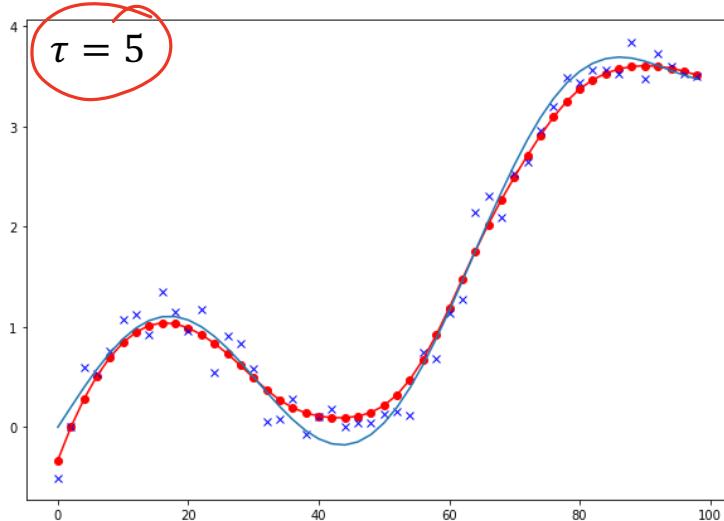
$$\text{Total Error} = \text{(Bias)}^2 + \text{Variance} + \text{Irreducible Error}$$

*Reducible error*

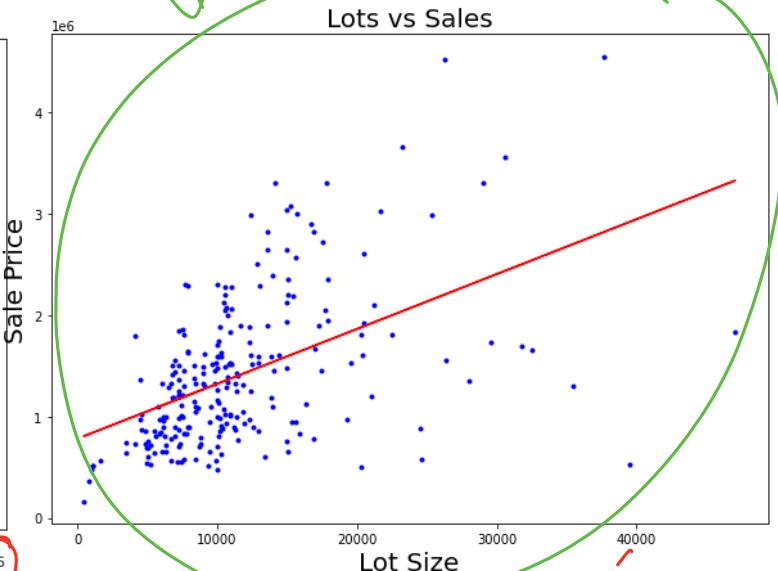
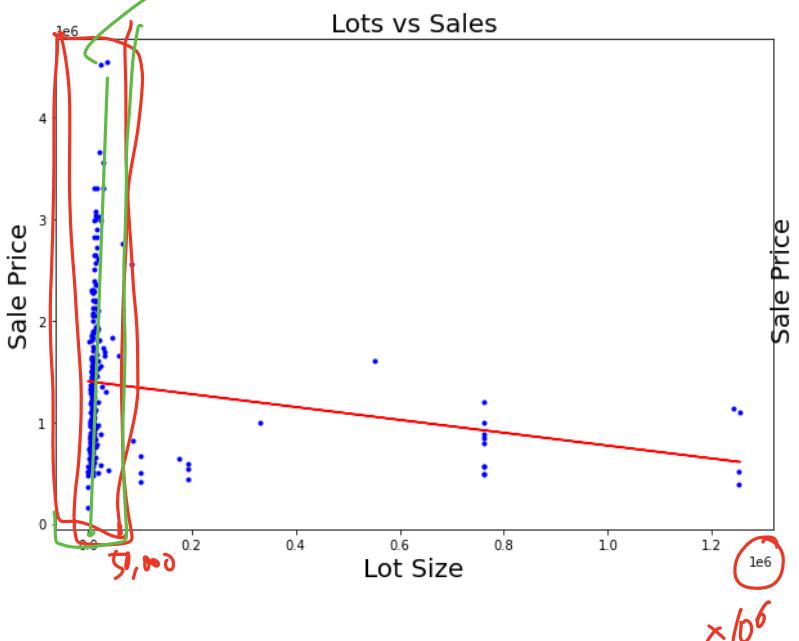
**Irreducible Error** is due to the noisiness of the data itself. The only way to reduce this part of the error is to clean up the data.



">// The best fit lies somewhere in between the extreme ends.



Remark: Clean up the data.

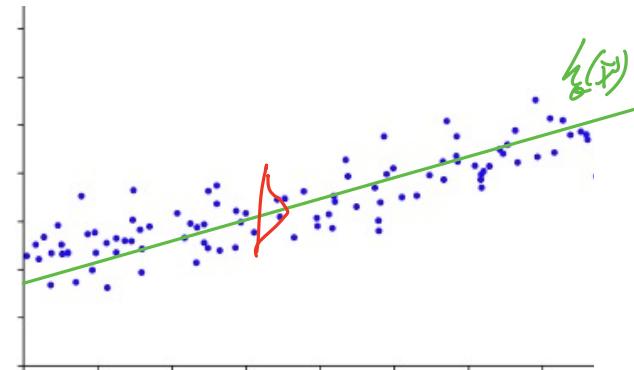


## ➤ Interpretation in Probability

**Data:**  $D = \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(n)}, y^{(n)})\}$

**Goal:** maximize the probability:  $P(y^{(i)} | \vec{x}^{(i)}; \theta)$

$$\underline{P(D_{\text{test}}|D_{\text{train}})}$$



Suppose the data follows linear model  $y = h(\vec{x}) + \epsilon$  with unmodeled error  $\epsilon$ .

This unmodeled error  $\epsilon$  is also called **irreducible error**. It is a random error term, which is independent from  $\vec{x}$  and has mean  $E(\epsilon)=0$ .

$$\epsilon \sim \text{Normal}(0, \sigma^2)$$

Even if it was possible to form a perfect estimate for  $h(\vec{x})$ , so that our estimated response took the form, our prediction would still have some error in it! The reason is that  $\epsilon$  does not come from the model. we cannot reduce the error introduced by  $\epsilon$ . The quantity  $\epsilon$  may also contain unmeasurable variation.

Our focus is on techniques for estimating  $h(\vec{x})$  with the aim of minimizing the reducible error comes from the model. Keep in mind that the irreducible error will always provide an upper bound on the accuracy of our prediction for  $y = h(\vec{x})$ . This bound is almost always unknown in practice.

## ➤ Linear Regression

$$(\vec{x}^{(i)}, y^{(i)})$$

peche

$$y = h(\vec{x}) + \epsilon = h_{\theta}(\vec{x}) + \epsilon$$

Suppose the data follows linear model  $y = \theta^T \vec{x} + \epsilon$  with unmodeled error  $\epsilon$ .

Suppose the unmodeled errors  $\epsilon^{(i)}$  are independent identical distribution(IID).

$$y^{(i)} = \theta^T \vec{x}^{(i)} + \epsilon^{(i)} \quad i=1, 2, \dots, n.$$

A more restrictive (but common) assumption: **Suppose** the unmodeled errors follow normal distribution,

$$\boxed{\epsilon \sim \text{Normal}(0, \sigma^2)}$$

- Given  $\vec{x}^{(i)}$  (constant)  
(Fixed)

$$\vec{y}^{(i)} | \vec{x}^{(i)} \sim \text{Normal}(\vec{\theta}^T \vec{x}^{(i)}, \sigma^2)$$

That is  $f(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \left(\frac{\epsilon^{(i)}}{\sigma}\right)^2\right)$

So,  $\underline{f(y^{(i)} | \vec{x}^{(i)}; \theta)} = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \left(\frac{y^{(i)} - \theta^T \vec{x}^{(i)}}{\sigma}\right)^2\right)$

➤ Interpretation in Probability

Maximize

$$P(\vec{y} | X)$$

Maximum Likelihood function:

$$\boxed{\text{Maximize}} \quad L(\vec{\theta}) := f(\vec{y} | X; \vec{\theta}) \underset{\text{iid}}{=} \prod_{i=1}^n f(y^{(i)} | \vec{x}^{(i)}; \vec{\theta})$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{y^{(i)} - \theta^T \vec{x}^{(i)}}{\sigma}\right)^2\right)$$

Find  $\vec{\theta}$  to

$$\boxed{\text{Maximize}} \quad \ln(L(\vec{\theta})) = -n \ln(\sqrt{2\pi}\sigma) - \sum_{i=1}^n \frac{1}{2} \left(\frac{y^{(i)} - \theta^T \vec{x}^{(i)}}{\sigma}\right)^2$$

$\iff$

$$\boxed{\text{Minimize}} \quad \sum_{i=1}^n (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2 = \text{RSS}(\vec{\theta})$$

Goal: Find  $\vec{\theta}^{\text{opt}} = \arg \max_{\vec{\theta}} \text{RSS}(\vec{\theta})$

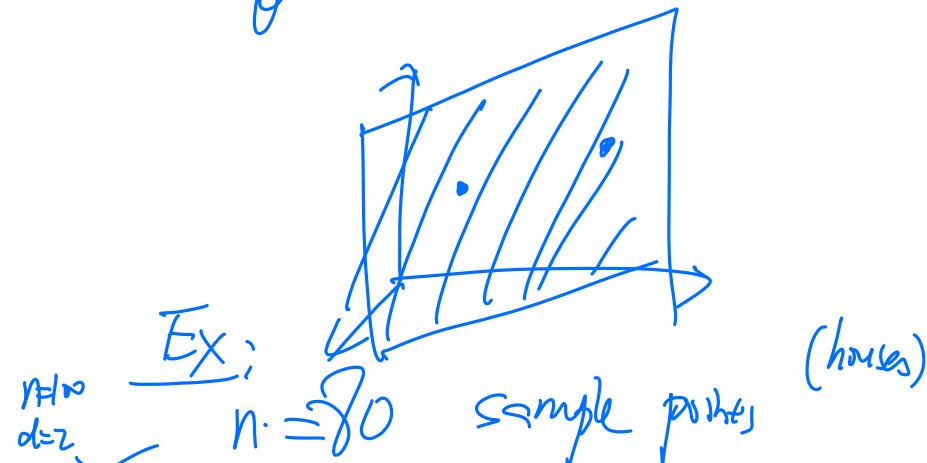
## Generalized linear model

and Exponential family

\* [Bishop] §2.4

\* [Murphy] §9.2, 9.3

$$= \underset{\vec{\theta}}{\arg \min} RSL(\vec{\theta})$$



data

$$\vec{X} \vec{\theta} = \vec{y}$$

$n \times d+1$   $d+1 \times 1$   $n \times 1$

$d=100$  features

first bedrooms  
last lot size

$x_1 = - -$   
 $\vdots$   
 $\vdots$

linear model

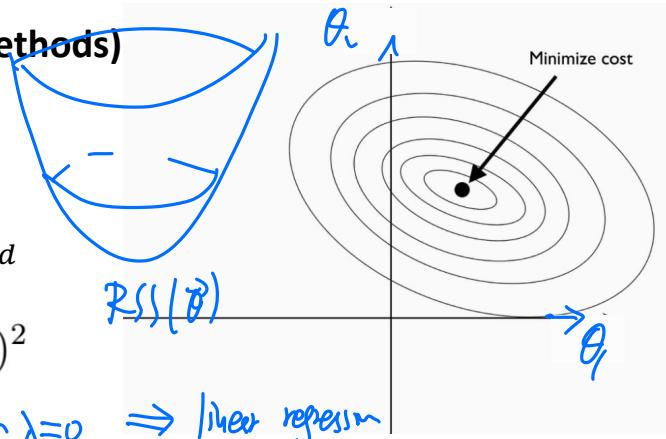
$$h_{\vec{\theta}}(\vec{x}) = \vec{\theta}^T \vec{x}$$

## ➤ Feature selection (Shrinkage/regularization methods)

Data:  $D = \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(n)}, y^{(n)})\}$

Model:  $h(\vec{x}) = \vec{\theta}^T \vec{x} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$

Square sum:  $\text{RSS}(\vec{\theta}) = \sum_{i=1}^n (h(\vec{x}^{(i)}) - y^{(i)})^2$



## ➤ Ridge Regression

Ridge

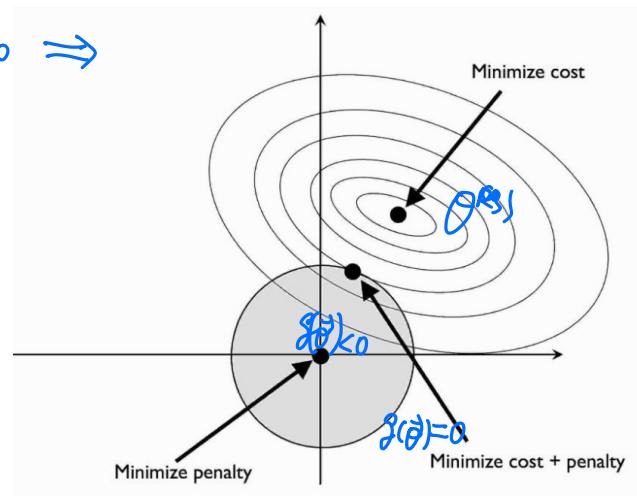
Ridge regression cost function:

$$J(\vec{\theta}) = \sum_{i=1}^n (y^{(i)} - h(\vec{x}^{(i)}))^2 + \lambda \sum_{j=1}^d \theta_j^2$$

cost ||| penalty

Lagrange multiplier

(KKT condition)



$$\vec{\theta}^{\text{ridge}} = \arg \min_{\vec{\theta}} J(\vec{\theta}) = \arg \min_{\vec{\theta}} \text{RSS}(\vec{\theta}) \quad \text{subject to} \quad g(\vec{\theta}) := \sum_{i=1}^d \theta_i^2 - t \leq 0$$

## ➤ Remarks on Ridge Regression

1. The first term measures goodness of fit, the smaller the better.
2. The second term is called shrinkage penalty, which shrinkage  $\theta_i$  towards to 0.
3. The shrinkage reduces variance (at the cost increased bias). Ridge regression works best in situations where the least squares estimates have high variance.
4. The intercept  $\theta_0$  is not penalized.
5. Ridge Regression is affected by the scale. (Least squares solution is unaffected by the scale.)
6. Ridge regression also has substantial computational advantages.

$$\underline{(x_i \theta_j)} = (c x_j) \underline{(\theta_j)}$$

$$J(\vec{\theta}) = (\vec{y} - X\vec{\theta})^T (\vec{y} - X\vec{\theta}) + \lambda \vec{\theta}^T \vec{\theta}$$

=====

$$\text{H.W.2} \quad \nabla J(\vec{\theta}) = 0$$

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

## Remark: (Standardization of feature variables/ Feature Scaling)

Before applying the Ridge/Lasso/ Elastic net regressions, we need to rescale an original variable to have equal range or variance.

### 1. Min-max scaling/normalization/ 0-1 scaling (Scikit-Learn: MinMaxScaler)

$$\frac{x_i - \min(x)}{\max(x) - \min(x)}$$

sklearn  
for each feature

### 2. Standardization (Scikit-Learn: StandardScaler)

$$\frac{x_i - \text{mean}(x)}{s(x)}$$

$$\bar{x}_i = \frac{x_i^{(1)} + \dots + x_i^{(n)}}{n}$$

$$s(x_i) = \sqrt{\sum_{j=1}^n \frac{(x_i^{(j)} - \bar{x}_i)^2}{n-1}}$$

where  $s(x_i)$  is the standard deviation of  $x$ .

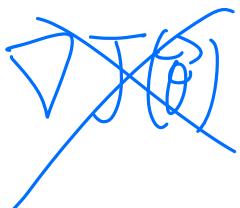
## ➤ Lasso Regression

Lasso regression cost function:

$$J(\vec{\theta}) = \text{RSS}(\vec{\theta}) + \lambda \sum_{j=1}^d |\theta_j|$$

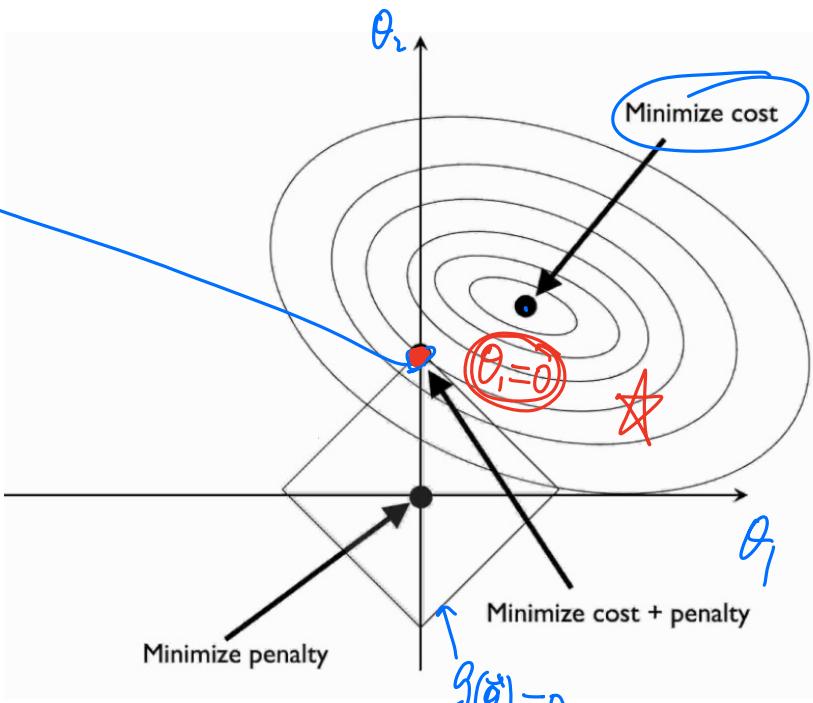
minimize

NO!



Goal:

$$\vec{\theta}^{(\text{lasso})} := \arg \min_{\vec{\theta}} J(\vec{\theta}) = \arg \min_{\vec{\theta}} \text{RSS}(\vec{\theta}) \text{ subject to } g(\vec{\theta}) = \sum_{j=1}^d |\theta_j| - t \leq 0$$



## • When $P > n$

$P = \# \text{ of features}$       Solve  $(X\vec{\theta} = \vec{y})$

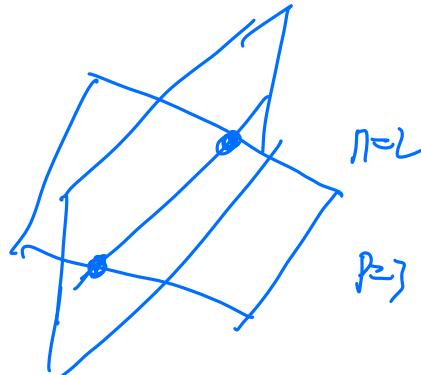
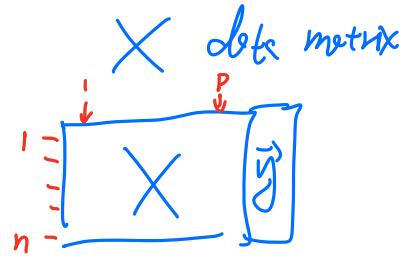
$\vec{x} \in \mathbb{R}^P$        $P \text{ large}$

- Solve  $X\vec{\theta} = \vec{y}$  is solvable.

• infinite many solu. for  $\vec{\theta}$ .

• infinite variance.

• zero bias (empty)



Goal Reduce variance, increase bias

Methods: ① Shrinking : ridge / LASSO.

② Subset selection      SVD.

③ dim reduction (PCA)

$$x_1 \dots x_p$$



$\hat{z}_i$  linear combination

$$\hat{z}_1, \hat{z}_2, \hat{z}_3 \dots$$

## ➤ Remarks on Lasso and Ridge Regressions

0. Lasso stands for Least Absolute Shrinkage and Selection Operator.
1. Lasso tends to completely eliminate the weights of the least important features. (It performs variable selection, and yields sparse models.) Hence, Lasso is more interpretable than ridge.
2. The lasso implicitly assumes that a number of the coefficients truly equal zero.
3. Ridge regression outperforms the lasso in terms of prediction error.
4. Both ridge and Lasso can improve over the traditional least squares by trade off variance with bias.
5. There are significant improvement when the variance of the least squares is large, mostly with small n and large d. = # of features
6. Lasso has feature selection, while ridge does not.
7. Use cross validation to determine which one has better prediction.
8. Ridge has closed form solution. Lasso generally does not have a closed form solution.

$$\vec{\theta} = \dots$$

## ➤ Elastic net Regression

Elastic net regression cost function:

$$J(\vec{\theta}) = \text{RSS}(\vec{\theta}) + \lambda \sum_{j=1}^d |\theta_j| + \eta \sum_{j=1}^d \theta_j^2$$

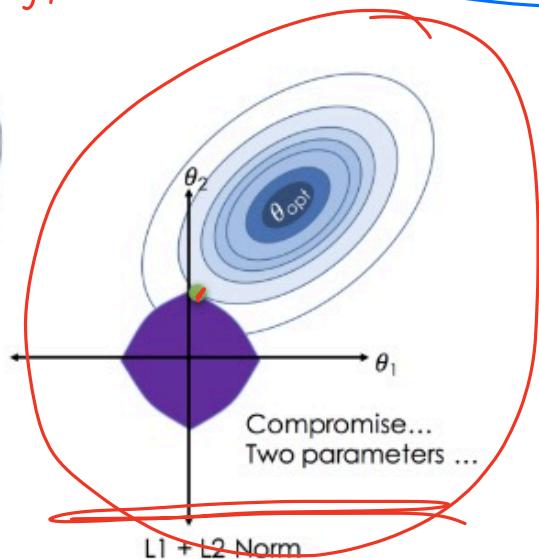
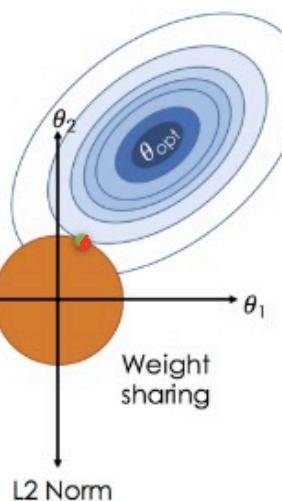
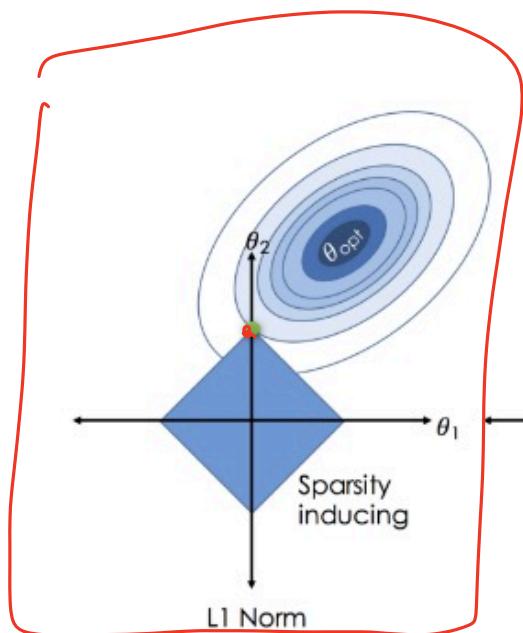
$\lambda \left( \sum_{j=1}^d |\theta_j| + \frac{1}{\lambda} \sum_{j=1}^d \theta_j^2 \right)$

$$\begin{matrix} L_1 \\ L_2 \end{matrix}$$

$L_p - \text{norm}$

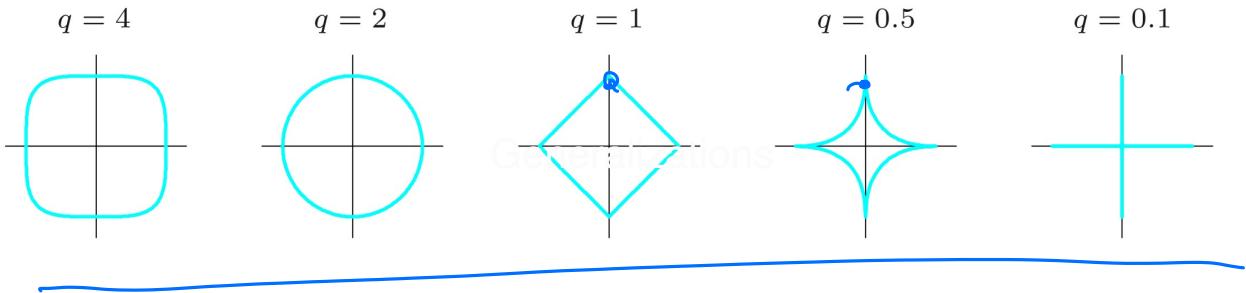
$$+ \lambda \sum_{j=1}^d |\theta_j|^p$$

$L_\infty - \text{norm}$  +  $\max \{|\theta_j| \cdot |\theta_j|\}$



**Generalizations:** For any positive number q.

$$J(\vec{\theta}) = \text{RSS}(\vec{\theta}) + \lambda \sum_{j=1}^d |\theta_j|^q$$



Generalizations

We introduced a few different statistical learning methods.  
Which one is the best approach?

George Box 1987: “All models are wrong, but some are useful.” (No Free Lunch Theorem)  
Hence, no one method dominates all others over all possible data sets. On a particular data set, one specific method may work best, but some other method may work better on a similar but different data set.

It is an important task to decide for any given set of data which method produces the best results. Selecting the best approach can be one of the most challenging parts of machine learning in practice. (Project?)