

$$1. EPE(f) = E[(Y - X^T \beta)^2]$$

$$= E[(Y - X^T \beta)^T (Y - X^T \beta)]$$

$$= E[Y^T Y - Y^T X^T \beta - \beta^T X Y + \beta^T X X^T \beta]$$

$$\frac{\partial EPE(f)}{\partial \beta} = 0 \Rightarrow E[-(Y^T X^T)^T - X Y + 2 X X^T \beta] = 0$$

$$E(X(Y - X^T \beta)) = 0$$

$$E(XY - X X^T \beta) = 0.$$

$$E(XY) = E(X X^T) \beta$$

$$\beta = E(X X^T)^{-1} E(XY)$$

2.

$$E[(f(\hat{x}_0) - \hat{y}_0)^2] = E[(f - \hat{y}_0)^2]$$

$$= E[f^2 - 2f\hat{y}_0 + \hat{y}_0^2]$$

$$= f^2 - 2fE(\hat{y}_0) + E(\hat{y}_0^2)$$

$$= f^2 - 2f\sigma^2 + 2fE(\hat{y}_0) + E(\hat{y}_0^2)$$

$$= (f - E(\hat{y}_0))^2 + E(\hat{y}_0) - E(\hat{y}_0)^2 + \sigma^2$$

(1)

Consider  $\vec{Y} = f(\vec{X}) + \epsilon = X^T \vec{\beta} + \epsilon$   $\vec{\beta}$  known.  $\epsilon_i \sim N(0, \sigma^2)$

$\hat{\vec{Y}} = X^T \hat{\vec{\beta}}$  fitted model.

$$r(\hat{\vec{\beta}}) = (X\hat{\vec{\beta}} - X\vec{\beta} - \vec{\epsilon})^T (X\hat{\vec{\beta}} - X\vec{\beta} - \vec{\epsilon})$$

$$= \hat{\vec{\beta}}^T X^T X \hat{\vec{\beta}} - \hat{\vec{\beta}}^T X^T X \vec{\beta} - \vec{\beta}^T X^T \epsilon - \vec{\beta}^T X^T X \hat{\vec{\beta}} + \vec{\beta}^T X^T X \vec{\beta} + \vec{\beta}^T X^T \epsilon - \epsilon^T X \hat{\vec{\beta}} + \epsilon^T X \vec{\beta} + \vec{\epsilon}^T \vec{\epsilon}$$

$$= (\hat{\vec{\beta}} - \vec{\beta})^T X^T X (\hat{\vec{\beta}} - \vec{\beta}) - 2\vec{\epsilon}^T X (\hat{\vec{\beta}} - \vec{\beta}) + \vec{\epsilon}^T \vec{\epsilon}$$

$$\frac{\partial r(\hat{\vec{\beta}})}{\partial \hat{\vec{\beta}}} = 0 \Rightarrow 2X^T X (\hat{\vec{\beta}} - \vec{\beta}) - 2X^T \vec{\epsilon} = 0$$

$$\Rightarrow \hat{\vec{\beta}} = \vec{\beta} + (X^T X)^{-1} X^T \vec{\epsilon}$$

$$\forall X_0 \quad \hat{y}_0 = \bar{X}_0^T \hat{\beta} = \bar{X}_0^T \beta + \bar{X}_0^T (X^T X)^{-1} X^T \epsilon \quad \textcircled{2}$$

Let  $l(\bar{X}_0) = X(X^T X)^{-1} X_0$  be the  $i$ th element of  $X(X^T X)^{-1} X_0$ .

By equation ①:

$$E[(f(\bar{X}_0) - \hat{y}_0)^2] = \underbrace{E[(f - E(\hat{y}_0))^2]}_{\text{bias}^2} + \underbrace{E[(\hat{y}_0 - E(\hat{y}_0))^2]}_{\text{variance}} + \sigma^2$$

$$\text{bias}^2 = E[(f - E(\hat{y}_0))^2] = \bar{X}_0^T \beta - \bar{X}_0^T \beta = 0. \quad \text{By } \textcircled{2}.$$

$$\begin{aligned} E(\hat{y}_0^2) - E(\hat{y}_0)^2 &= \text{Var}(\hat{y}_0) = \text{Var}(\bar{X}_0^T \hat{\beta} + \bar{X}_0^T (X^T X)^{-1} X^T \epsilon) \\ &= \text{Var}(\bar{X}_0^T (X^T X)^{-1} X^T \epsilon) \\ &= \text{Var}\left(\sum_{i=1}^N l_i(\bar{X}_0) \epsilon_i\right) \\ &= \sum_{i=1}^N l_i^2(\bar{X}_0) \text{Var}(\epsilon_i) \\ &= E_T(\bar{X}_0^T (X^T X)^{-1} \bar{X}_0) \sigma^2. \end{aligned}$$

$$\Rightarrow: EPE(X_0) = E_T \bar{X}_0^T (X^T X)^{-1} \bar{X}_0 \sigma^2 + 0^2.$$

3.  
a)  $X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$   $X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}$   $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

In class, we've proved that  $\hat{\beta} = (X^T X)^{-1} X^T y$ .

From the form, we know

In linear regression:  $l_i(x_0; X) = [1 \ x_0] (X^T X)^{-1} \begin{bmatrix} 1 \\ x_i \end{bmatrix}$  for  $k \leq n$ .

where  $X^T X = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$

When  $\det(X^T X) > 0$  We can compute the inverse.

In k-nearest neighbors:

$$l_i(x_0; X) = \begin{cases} \frac{1}{k} & \text{if } x_i \text{ is one of the nearest } k \text{ points.} \\ 0 & \text{otherwise.} \end{cases}$$

b).  $X$  is fixed  $y$  varies.

$$\begin{aligned} E_{y|X}((f(x_0) - \hat{f}(x_0))^2) &= f(x_0)^2 - 2f(x_0) E_{y|X}(f(\hat{x}_0)) + E_{y|X}(f(\hat{x}_0)^2) \\ &= (f(x_0) - E_{y|X}(f(\hat{x}_0)))^2 + E_{y|X}(f(\hat{x}_0)^2) - E_{y|X}(f(\hat{x}_0))^2 \\ &= \text{bias}^2 + \text{Var}(f(\hat{x}_0)). \end{aligned}$$

```

In [39]: from numpy.linalg import inv
from sklearn.metrics import r2_score
from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import LinearRegression
import numpy as np
from sklearn.metrics import mean_absolute_error, mean_squared_error
import matplotlib.pyplot as plt
%matplotlib inline
for i in range(3):
    print("In the ", i+1, " trial")
    x = np.random.uniform(0,1,(100,1))*5-2##### o-1??
    mu, sigma = 0, 1 # mean and standard deviation
    s = np.random.normal(mu, sigma, (100,1))
    Y = 0.5*x*x+0.5*x+1+s
    plt.plot(x,Y,'bx')
    xsort = np.sort(x)
    line = np.array(np.sort(x,axis = None)).reshape(100,1)
    for ii in [1,2,7,50]:
        poly = PolynomialFeatures(ii, include_bias=False)
        tran_X = poly.fit_transform(xsort)
        reg = LinearRegression().fit(tran_X, Y)
        line_poly = poly.transform(line)
        y_hat = reg.predict(line_poly)
        mae = mean_absolute_error(Y, y_hat)
        print("p= ", ii, "mean_absolute_error is", mae)
        plt.plot(line, y_hat)
    plt.show()

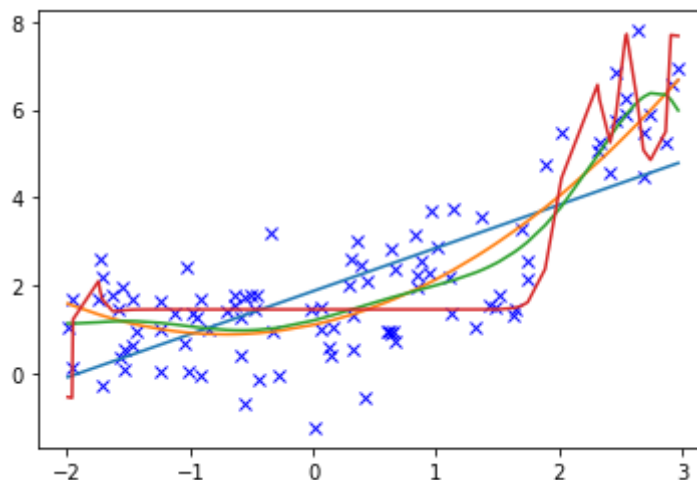
```

In the 1 trial

```

p= 1 mean_absolute_error is 1.7674505564609846
p= 2 mean_absolute_error is 1.803352402078023
p= 7 mean_absolute_error is 1.811035948348298
p= 50 mean_absolute_error is 1.8273915868630255

```

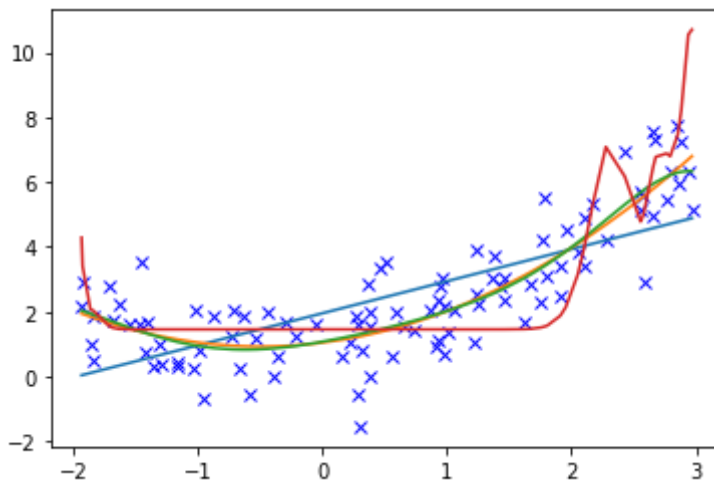


In the 2 trial

```

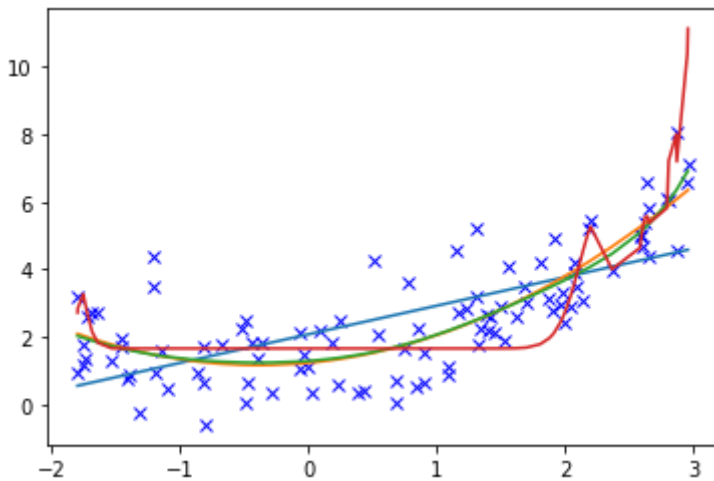
p= 1 mean_absolute_error is 1.948030411837093
p= 2 mean_absolute_error is 2.0727259430897997
p= 7 mean_absolute_error is 2.07665520059751
p= 50 mean_absolute_error is 2.0463238109468147

```



In the 3 trial

```
p= 1 mean_absolute_error is 1.7285544609961285
p= 2 mean_absolute_error is 1.8497399443649098
p= 7 mean_absolute_error is 1.852059503268694
p= 50 mean_absolute_error is 1.934517780639168
```



b) From the graph, we see both  $p=2$  and  $p=3$  give us the best fitting result. However,  $p=7$  has downward trend in graph 1, so I pick  $p = 2$  as the best. If I didn't know the true regression function and guess the model from the resulting graph, I'll pick  $p = 1$  as the best because  $p=1$  has the lowest absolute error compared to the last 3 in any trial.

In [ ]: