

---

---

---

---

---



Q-1 (a)

$$E[(Y - X'b)^2] = \langle Y - X'b, Y - X'b \rangle$$

$$= \langle U + X'\beta_0 - X'b, U + X'\beta_0 - X'b \rangle$$

$$= \langle U + X'(\beta_0 - b), U + X'(\beta_0 - b) \rangle$$

$$= \langle U, U \rangle + \langle X'(\beta_0 - b), U \rangle + \langle X'(\beta_0 - b), U \rangle + \langle X'(\beta_0 - b), X'(\beta_0 - b) \rangle$$

$$= E[U^2] + 2 \langle X'(\beta_0 - b), U \rangle + \langle X'(\beta_0 - b), X'(\beta_0 - b) \rangle$$

$$= E[U^2] + 2 E[(X'(\beta_0 - b)U)] + E[(X'(\beta_0 - b))^2]$$

$$= E[U^2] + 2 E[(\beta_0 - b)' X U] + E[X'(\beta_0 - b) X'(\beta_0 - b)]$$

When  $X'(\beta_0 - b)U$  is scalar, so its transpose is equal to itself.

$$(AB)' = B'A'$$

Similarly,  $(X'(\beta_0 - b))' = (\beta_0 - b)'X$

$$E[(Y - X'b)] = E[U^2] + 2(\beta_0 - b)'E[XU] + (\beta_0 - b)'E[XX'](\beta_0 - b)$$



(b)

We can write quadratic form,

$$\begin{aligned} (\beta_0 - b)'E[XX'](\beta_0 - b) &= E[(\beta_0 - b)'XX'(\beta_0 - b)] \\ &= E[(\beta_0 - b)'X]^2 \end{aligned}$$

It is given  $E[(\beta_0 - b)'X]^2 \geq 0$  & is strictly positive if  $\beta_0 \neq b$ . So,

$$\begin{aligned}
 E[(Y - X'b)^2] &= E[U^2] + 2(\beta_0 - b)'E[XU] + E[(\beta_0 - b)'X]^2 \\
 &= E[U^2] + 0 + \underbrace{E[(\beta_0 - b)'X]^2}_{\geq 0} \\
 &\geq E[U^2]
 \end{aligned}$$

As we know,  $E[(\beta_0 - b)'X]^2 > 0$  if  $\beta_0 \neq b$ . So, strict inequality for  $\beta_0 \neq b$ .

As we see squared error of a linear predictor of  $y$  is  $(Y - X'b)^2$  and its expectation is  $E[(Y - X'b)^2]$ . Also, this expectation is minimized when  $b = \beta_0$  i.e. when we use population linear predictor.

(c)

$$\begin{aligned} V(Y) &= E[(Y - E[Y])^2] \\ &= E[(E^*[Y|X] + U - E[Y])^2] \\ &= E[(E^*[Y|X] - E[Y] + U)^2] \\ &= E\left\{ (E^*[Y|X] - E[Y])^2 + 2(E^*[Y|X] - E[Y])(U) + U^2 \right\} \end{aligned}$$

Since,  $E^*[Y|X]$  is projection, so projection error is orthogonal to  $E^*[Y|X]$ . Also,  $E[U] = 0$ , so cross term is 0.

$$\text{Also, } E[Y] = E[E^*[Y|X]]$$

$$\text{Var}(Y) = E[(E^*[Y|X] - E[E^*[Y|X]])^2] + E[U^2]$$

$$\text{Var}(Y) = \text{Var}(E^*[Y|X]) + E[U^2]$$

but  $E[\epsilon^2] = \text{Var}(Y - E^*[Y|X])$ . So,

$$\text{Var}(Y) = \text{Var}(E^*[Y|X]) + \text{Var}(Y - E^*[Y|X])$$

So, total variance of  $Y$  variable due to linear component/linear projection plus the part not explained by linear projection.

(d)

Since  $X_K$  is scalar RV,

$$E^*[Y|X_K] = \alpha_K + \beta_K X_K$$

Since  $X_K, X_Q$  are uncorrelated &  $E^*[Y|X_K]$  is linear function of  $X_K$ ,

$$\text{Cov}(E^*[Y|X_K], X_Q) = 0$$

but since two are uncorrelated, projection  
of this onto  $X_1$ ,

$$E[E[Y|X_k]|X_1] = E[E[Y|X_k]] = E[Y]$$

$$E[UX_1] = E[YX_1] - \sum_{k=1}^K E[E[Y|X_k]|X_1] + (K-1)E[Y]E[X_1]$$

Since  $E[Y|X_k]$  is lin,

$$\begin{aligned} E[E[Y|X_k]|X_1] &= E[E[Y|X_k]]E[X_1] \\ &= E[Y]E[X_1] \end{aligned}$$

for  $k=1$ ,

$$E[E[Y|X_k]|X_1] = E[YX_1]$$

$\Rightarrow$

$$\begin{aligned} E[UX_1] &= E[YX_1] - E[YX_1] - (K-1)E[Y]E[X_1] \\ &\quad + (K-1)E[Y]E[X_1] \end{aligned}$$

$$E[UX] = 0 + 0 = 0 \quad \forall$$

Note that when  $X_K$  is one regressed  
 then

$$E^*[Y | X_K] = \alpha_K + \beta_K X_K$$

$$\text{where } \alpha_K = E[Y] - \beta_K E[X_K]$$

$$\beta_K = \frac{\text{Cov}(Y, X_K)}{V(X_K)}$$

So,

$$E^*[Y | X_1, \dots, X_N] = \sum_{k=1}^N E^*[Y | X_K] - (K-1) E[Y]$$

$$= E[Y] + \sum_{k=1}^N (E^*[Y | X_K] - E[Y])$$

$$= E[Y] + \sum_{k=1}^N (\alpha_K + \beta_K X_K - E[Y])$$



$$= E[Y] + \sum_{k=1}^N (E[Y] - \beta_k E[X_k] + \beta_k X_k - E[Y])$$

$$= E[Y] + \sum_{k=1}^K \beta_k (X_k - E[X_k])$$

$$= E[Y] + \sum_{k=1}^K \frac{\text{Cov}(Y, X_k)}{\text{Var}(X_k)} (X_k - E[X_k])$$

Now,

$$\text{Var}(E[Y|X_1, \dots, X_K]) = \text{Var}\left(\sum_{k=1}^K \frac{C(Y, X_k)}{\text{Var}(X_k)} X_k\right)$$

as variance of expectation is 0.

$$= \sum_{k=1}^K \frac{C(Y, X_k)^2}{\text{Var}(X_k)} \text{Var}(X_k) +$$

$$+ \sum_{k \neq i} \frac{\text{Cov}(X_k, X_i) \text{Cov}(Y, X_i)}{\text{Var}(X_k) \text{Var}(X_i)} C(X_k, Y_i)$$

$$C(X_k, X_i) = 0, \text{ so,}$$

$$\text{Var}(E^\circ[Y|X_1, \dots, X_N]) = \sum_{k=1}^N \frac{\text{Cov}(Y, X_k)^2}{\text{Var}(X_k)}$$

$$\text{Now, } 1 - \frac{\text{Var}(U)}{\text{Var}(Y)} = \frac{\text{Var}(Y) - \text{Var}(U)}{\text{Var}(Y)}$$

we showed in part (c),

$$\text{Var}(Y) - \text{Var}(U) = \text{Var}(E^\circ[Y|X])$$

$$1 - \frac{\text{Var}(U)}{\text{Var}(Y)} = \frac{\text{Var}(E^\circ[Y|X_1, \dots, X_N])}{\text{Var}(Y)}$$

$$= \sum_{k=1}^N \frac{\text{Cov}(Y, X_k)^2}{\text{Var}(Y) \text{Var}(X_k)}$$

$$= \sum_{k=1}^N \rho_k^2$$



$$Y_i \sim N(\mu, \sigma_i^2) \quad \text{Q-2 (a)}$$

$$\text{choose } c = \begin{pmatrix} 1/N \\ 1/N \\ \vdots \\ 1/N \end{pmatrix} \quad ( \quad )$$

$$\text{then } c'Y = \frac{1}{N}Y_1 + \frac{1}{N}Y_2 + \dots + \frac{1}{N}Y_N = \bar{Y}$$

so, there exists such a  $c$ .

(b)

$$MSE(\hat{\mu}) = \text{Bias}^2 + \text{Var}(\hat{\mu})$$

$$E[\hat{\mu}] = c' E[Y] = c' \mu \mathbf{1}_N$$

$$\text{Var}(\hat{\mu}) = c' \text{Var}(Y) c = c' \text{diag}(\sigma_1^2, \dots, \sigma_N^2) c$$

$$\text{Bias}(\hat{\mu}) = \mu c' \mathbf{1}_N - \mu \mathbf{1}_N$$

$$MSE(\hat{c}) = c' \text{diag}(\sigma_1^2, \dots, \sigma_N^2) c + (\mathcal{H}'c - \mathcal{H})$$

Taking matrix derivative,

$$\frac{\partial MSE(\hat{c})}{\partial c} = 2(\text{diag}(\sigma_1^2, \dots, \sigma_N^2))c + 2\mathcal{H}'c - \mathcal{H}$$

$$\text{diag}(\sigma_1^2, \dots, \sigma_N^2)c + \mathcal{H}'c - \mathcal{H} = 0 \quad = 0$$

$$[\text{diag}(\sigma_1^2, \dots, \sigma_N^2) + \mathcal{H}'\mathcal{H}]c = \mathcal{H}'\mathcal{H}$$

$$c = [\text{diag}(\sigma_1^2, \dots, \sigma_N^2) + \mathcal{H}'\mathcal{H}]^{-1} \mathcal{H}'\mathcal{H}$$

(C)

$$\sigma_i^2 = \sigma^2.$$

Let us call,  $D = \text{diag}(\sigma^2, \dots, \sigma^2)$

$$\mathbf{1}'_N D^{-1} \mathbf{1}_N = \sum_{i=1}^N \frac{1}{\sigma^2}$$

$$C = (D + \mathcal{H}^2 \mathbf{1}_N \mathbf{1}'_N)^{-1} \mathcal{H}^2 \mathbf{1}_N$$

$$(D + \mathcal{H}^2 \mathbf{1}_N \mathbf{1}'_N)^{-1} = D^{-1} - \frac{\mathcal{H}^2 D^{-1} \mathbf{1}_N \mathbf{1}'_N D^{-1}}{1 + \mathcal{H}^2 \mathbf{1}'_N D^{-1} \mathbf{1}_N}$$

(d)

$$\text{Bias}(\hat{\mu}) = \left[ \frac{\mu^2}{\frac{\sigma^2}{N} + \mu^2} - 1 \right] \mu$$

$$\text{Var}(\hat{\mu}) = \left( \frac{\mu^2}{\frac{\sigma^2}{N} + \mu^2} \right)^2 \cdot \frac{\sigma^2}{N}$$

When  $N \rightarrow \infty$ ,  $\text{Bias}(\hat{\mu}) \rightarrow 0$  as well as  $\text{Var}(\hat{\mu}) \rightarrow 0$ , so it converges in probability to  $\mu$ .

$$\frac{\mu^2 \bar{D}' (N \bar{D}' \bar{D})^{-1}}{1 + \mu^2 \sum_{i=1}^N 1/\sigma^2} =$$

$$\text{Since, } \hat{\mu} = \frac{\mu^2}{\sigma^2/N + \mu^2} \bar{Y}$$

$$\text{when } N \rightarrow \infty, \bar{Y} \rightarrow \mu, \sigma^2/N \rightarrow 0$$

$$\text{So, } \hat{\mu} \xrightarrow{P} \frac{\mu^2}{\mu^2} \cdot \mu = \mu$$

(e)

Since that estimator relied on true mean  $\mu$ , here we can replace true mean  $\mu$  also with  $\bar{Y}$ . Thus,

We can write

$$\hat{\mu} = \frac{\bar{Y}^2}{\frac{\sigma^2}{N} + \bar{Y}^2} \bar{Y}$$

if we know  $\frac{\sigma^2}{N}$  is very small,

$$\text{th } \hat{\mu} = \frac{\bar{Y}^2 - \frac{\sigma^2}{N}}{\bar{Y}^2} \bar{Y}$$

will be similar (cannot think of a good reason to do this).

The MSE of this will be lower as we showed in class MLE is inadmissible in this type of estimators and above shrinkage estimator has uniformly lower maximum MSE.



(F)

Mean is still pretty good estimator  
and especially in large sample all of  
these have similar performance.

Also for mean you don't need to know  
any extra information and it is  
much easier to calculate.

# Problem 3 PS 2

Muhammad Bashir

```
In [1]: # Load Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

In [3]: # set paths to working directory
path = '/Users/muhammadbashir/GitHub/MuhammadCourses/Ec240a/Problem Sets'
# load RPS_calorie_data.out data and read only columns Y0tc and X0te.
calories = pd.read_csv(path + '/RPS_calorie_data.out',usecols=['Y0tc','X0te'])
calories.head()
```

Out[3]:

	Y0tc	X0te
0	9.061790	9.917698
1	7.870896	8.276105
2	8.262976	8.036979
3	9.057076	9.212327
4	8.953918	9.826068

Let ( Y ) denote log calories and ( X ) denote log expenditure. Assume that  $[ m(x) = \mathbb{E}[Y \mid X = x] = \sum_{k=1}^K \alpha_k g_k(x), ]$

where (  $g_k(x) = x^{k-1}$  ).

Using the power series basis described above and the Gram-Schmidt algorithm construct a new basis that is orthogonal to the design points (set (  $K = 12$  )). Let (  $W_i$  ) denote the (  $K \times 1$  ) vector of orthonormal basis functions for the ( i )-th household. Compute the least squares fit (  $m(X_i) = W_i^0 \hat{\theta}$  ) with  $[ \hat{\theta} = \left( \sum_{i=1}^N W_i W_i^0 \right)^{-1} \times \left( \sum_{i=1}^N W_i Y_i \right). ]$  Plot this function onto a scatter of the unsmoothed data.

```
In [4]: # First create K =12 g functions/polynomials, where g_k(x) = x^k-1 for the column X0te.
K = 12
```

```
for k in range(1,K+1):
    calories['X0te_g'+str(k)] = calories['X0te']**k
```

Orthogonalize g functions using Gram-Schmidt algorithm

```
In [5]: def gram_schmidt_custom_inner_product(X):
        """
        Orthonormalizes the columns of matrix X using the Gram-Schmidt process
        with a custom inner product defined as the mean of the element-wise product.

        Parameters:
            X (numpy.ndarray): Input matrix with shape (m, n), where m >= n.

        Returns:
            numpy.ndarray: Matrix with orthonormalized columns.
        """
        def custom_inner_product(u, v):
            """Custom inner product: mean of element-wise product."""
            return np.mean(u * v)

        # Get the number of rows (m) and columns (n) of X
        m, n = X.shape

        # Initialize an empty matrix to store orthonormal vectors
        Q = np.zeros((m, n))

        for i in range(n):
            # Start with the current column vector
            v = X[:, i]

            # Subtract projections of v onto all previous Q columns
            for j in range(i):
                projection = custom_inner_product(Q[:, j], v) * Q[:, j]
                v -= projection

            # Normalize the vector using the custom norm
            v_norm = np.sqrt(custom_inner_product(v, v))
            if v_norm > 1e-10: # Check to avoid division by zero
                Q[:, i] = v / v_norm
            else:
                raise ValueError(f"Column {i} is linearly dependent on previous columns.")
```

```
return Q
```

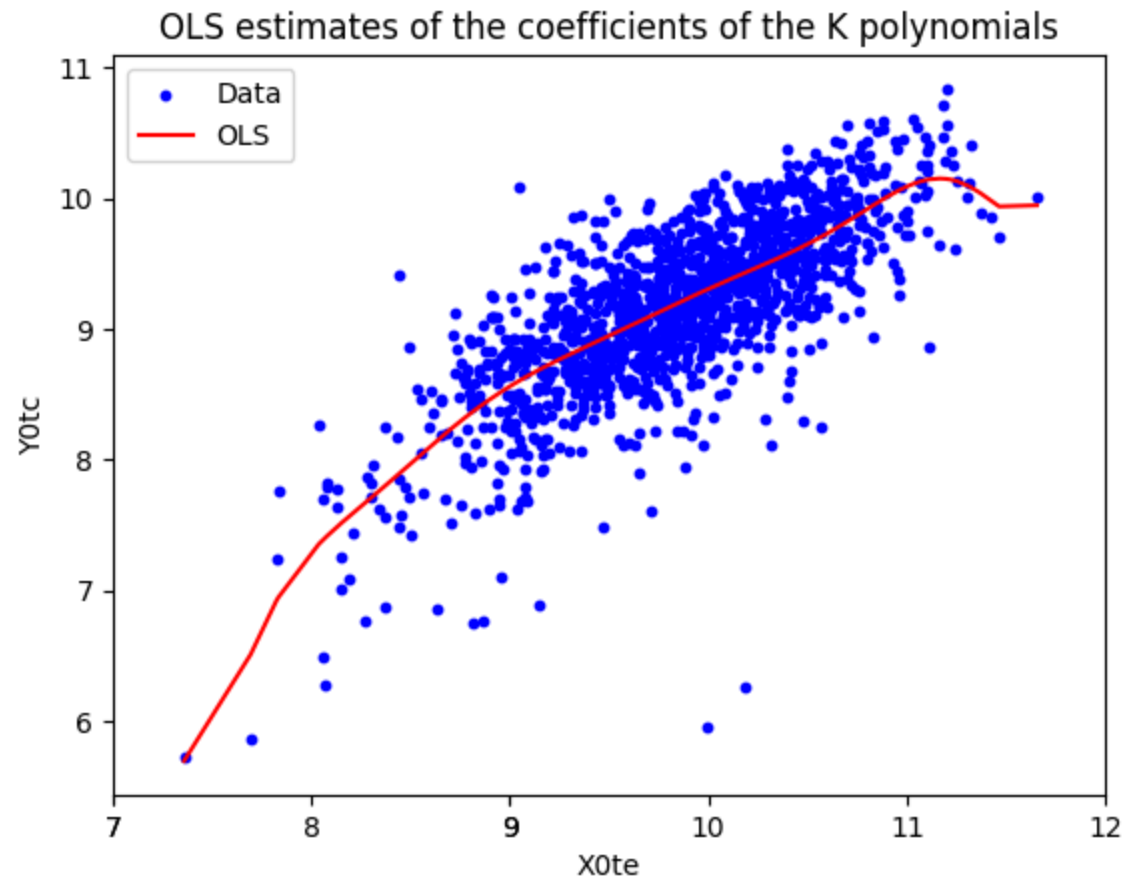
```
# Apply the Gram-Schmidt process to the columns of X0te  
X = calories[['X0te_g' + str(k) for k in range(1, K+1)]].values  
Q = gram_schmidt_custom_inner_product(X)  
Q.shape
```

Out[5]: (1358, 12)

Compute least square fit of this orthonormal basis

```
In [6]: # Compute the OLS estimates of the coefficients of the K polynomials i.,e regress Y0tc on the K polynomials fk.  
# First create a list of the columns of the data frame that contain the K polynomials.  
# X = calories[['X0te_f'+str(k) for k in range(1,13)]]  
# X = np.array(X)  
X= Q  
Y = calories['Y0tc']  
Y = np.array(Y)  
  
# Compute the OLS estimates of the coefficients of the K polynomials i.,e regress Y0tc on the K polynomials fk.  
beta = (np.linalg.inv(X.T @ X) @ X.T) @ Y  
# predict mean by multiplying the estimated coefficients by the polynomials.  
calories['Y0tc_hat'] = X @ beta
```

```
In [7]: # Plot this function onto a scatter of the unsmoothed data.  
plt.scatter(calories['X0te'], calories['Y0tc'], label='Data', color='blue', s=10)  
sorted_indices = np.argsort(calories['X0te'])  
plt.plot(calories['X0te'].iloc[sorted_indices], calories['Y0tc_hat'].iloc[sorted_indices], label='OLS', color='red')  
plt.xlabel('X0te')  
plt.ylabel('Y0tc')  
plt.legend()  
plt.title('OLS estimates of the coefficients of the K polynomials')  
  
# Add 10 equally spaced xticks  
xticks = np.linspace(np.floor(calories['X0te'].min()), np.ceil(calories['X0te'].max()), 8)  
xticks = xticks.astype(int)  
plt.xticks(xticks)  
  
plt.show()
```



3. Now use the shrinkage estimator of Efromovich (1999) as described in lecture to estimate  $m(X_i)$ . Plot this function onto a scatter of the unsmoothed data. Comment on your findings.

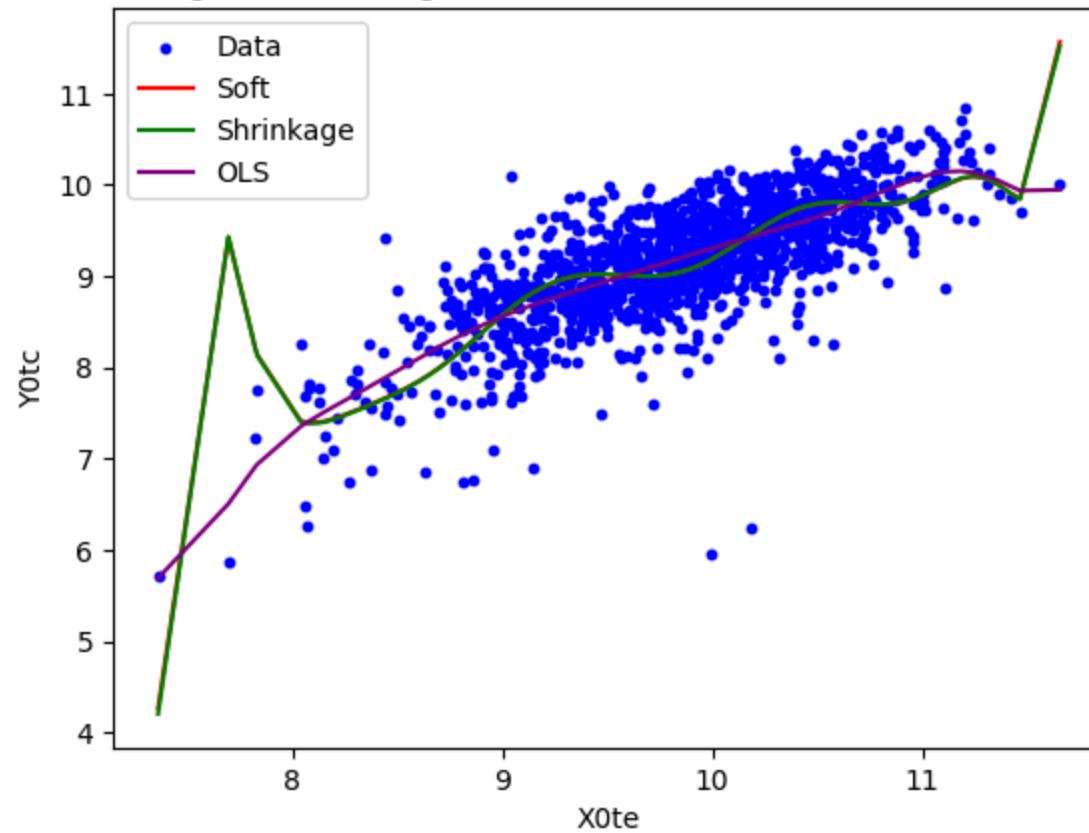
```
In [8]: # Define Shrinkage constant c as  $1 - (\sigma^2 / (N * \beta^2))$ 
sigma2 = np.var(Y - X @ beta)
N = len(Y)
c = np.zeros(K)
theta = np.zeros(K)
for i in range(K):
    c[i] = 1 - (sigma2 / N) / beta[i]**2
    theta[i] = c[i] * beta[i]
calories['Y0tc_hat_shrink'] = np.dot(X, theta)
```

## Soft Threshold Estimator

```
In [9]: # First let us choose lambda by minimizing  $(K/N) \cdot \sigma^2 - 2(\sigma^2/N) \sum_{l=1}^K (1(\text{abs}(\theta_k) \leq \lambda) + \sum_{l=1}^K (\theta_k^2, \lambda^2))$ 
# Define the function that computes the value of the objective function for a given value of lambda.
def objective_function(lmbda, theta, sigma2, K, N):
    return (K/N)*sigma2 - 2*(sigma2/N)*np.sum(np.abs(theta) <= lmbda) + np.sum(np.minimum(theta**2, lmbda**2))
# minimize the objective function over lambda
from scipy.optimize import minimize_scalar
result = minimize_scalar(objective_function, args=(theta, sigma2, K, N))
lmbda = result.x
# Soft thresholding
theta_soft = np.sign(theta) * np.maximum(np.abs(theta) - lmbda, 0)
calories['Y0tc_hat_soft'] = X @ theta_soft

In [10]: # Plot the soft thresholding estimates of the coefficients of the K polynomials
plt.scatter(calories['X0te'], calories['Y0tc'], label='Data', color='blue', s=10)
sorted_indices = np.argsort(calories['X0te'])
plt.plot(calories['X0te'].iloc[sorted_indices], calories['Y0tc_hat_soft'].iloc[sorted_indices], label='Soft', color='red')
plt.plot(calories['X0te'].iloc[sorted_indices], calories['Y0tc_hat_shrink'].iloc[sorted_indices], label='Shrinkage', color='green')
plt.plot(calories['X0te'].iloc[sorted_indices], calories['Y0tc_hat'].iloc[sorted_indices], label='OLS', color='purple')
plt.xlabel('X0te')
plt.ylabel('Y0tc')
plt.legend()
plt.title('Soft thresholding and Shrinkage estimates of the coefficients of the K polynomials')
plt.show()
```

Soft thresholding and Shrinkage estimates of the coefficients of the K polynomials



We notice that OLS has overfitting to some extent as it is too smooth. Both Soft threshold and shrinkage estimators are same but different from OLS. These estimators punish overfitting and hence look different from OLS. Notice that in the middle, red/green line is more flexible.