

- 1) a) In the linear regression case, the estimator of f at x_0 is given by
- $$\hat{f}(x_0) = x_0^T \hat{\beta}; \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

We can then rewrite $\hat{f}(x_0)$ as:

$$\hat{f}(x_0) = \sum_{i=1}^N (x_0^T (X^T X)^{-1} X^T)_i y_i.$$

Therefore,

$$\underline{\ell_i(x_0; X) = (x_0^T (X^T X)^{-1} X^T)_i.}$$

In K-Nearest-Neighbors, we see that

$$\hat{f}(x_0) = \frac{1}{K} \sum_{\{x_i \in N_K(x_0)\}} y_i = \sum_{i=1}^N y_i \cdot \frac{1}{K} \mathbb{1}_{\{x_i \in N_K(x_0)\}},$$

and clearly,

$$\underline{\ell_i(x_0; X) = \frac{1}{K} \cdot \mathbb{1}_{\{x_i \in N_K(x_0)\}}}$$

$$\begin{aligned} \text{b) } \mathbb{E}_{y|X} [f(x_0) - \hat{f}(x_0)]^2 &= \mathbb{E}_{y|X} [\hat{f}(x_0) - \mathbb{E}_{y|X}(\hat{f}(x_0))]^2 - \mathbb{E}_{y|X} [\mathbb{E}_{y|X}[\hat{f}(x_0)] - f(x_0)]^2 \\ &= \underline{\text{Var}_{y|X}(\hat{f}(x_0)) + \text{Bias}_{y|X}^2(\hat{f}(x_0))} \end{aligned}$$

$$\begin{aligned} \text{c) } \mathbb{E}_{y,X} [f(x_0) - \hat{f}(x_0)]^2 &= \mathbb{E}_{y,X} [\hat{f}(x_0) - \mathbb{E}_{y,X}(\hat{f}(x_0))]^2 - \mathbb{E}_{y,X} [\mathbb{E}_{y,X}[\hat{f}(x_0)] - f(x_0)]^2 \\ &= \underline{\text{Var}_{y,X}(\hat{f}(x_0)) + \text{Bias}^2(\hat{f}(x_0))} \end{aligned}$$

- d) From b and c, we see the previous variance-bias equations. Comparing them, we see:

$$\text{Var}_{y|X}(\hat{f}(x_0)) + \text{Bias}_{y|X}^2(\hat{f}(x_0)) = \text{Var}_{y,X}(\hat{f}(x_0)) + \text{Bias}_{y,X}^2(\hat{f}(x_0))$$

Taking the expectation w.r.t. X :

$$\begin{aligned} \mathbb{E}_X [\text{Var}_{y|X}(\hat{f}(x_0))] + \mathbb{E}_X [\text{Bias}_{y|X}^2(\hat{f}(x_0))] &= \mathbb{E} [\text{Var}_{y,X}(\hat{f}(x_0)) + \text{Bias}_{y,X}^2(\hat{f}(x_0))] \\ \Rightarrow \mathbb{E}_X [\text{Bias}_{y|X}^2(\hat{f}(x_0))] - \text{Bias}^2(\hat{f}(x_0)) &= \text{Var}(\hat{f}(x_0)) - \mathbb{E}_X [\text{Var}_{y|X}(\hat{f}(x_0))] \\ &= \text{Var}_{y|X}(\mathbb{E}_X[\hat{f}(x_0)]) \\ &\geq 0. \end{aligned}$$

This suggests that the expectation of the conditional squared bias is greater than or equal to the squared bias of $\hat{f}(x_0)$, and the difference is equal to $\text{Var}_{y|X}(\mathbb{E}_X[\hat{f}(x_0)])$.

2) let g be the natural cubic spline with knots @ x_i , and \tilde{g} be any other differentiable function on $[a, b]$ that interpolates the N pairs.

let $h(x) = \tilde{g}(x) - g(x)$. We show that $\int_a^b g''(x) h''(x) dx = 0$

$$\begin{aligned} \int_a^b g''(x) h''(x) dx &= \int_a^b g''(x) (\tilde{g}''(x) - g''(x)) dx = g''(x) (\tilde{g}'(x) - g'(x)) \Big|_a^b \\ &\quad - \int_a^b g'''(x) (\tilde{g}'(x) - g'(x)) dx \\ &= \underbrace{g''(b) (\tilde{g}'(b) - g'(b)) - g''(a) (\tilde{g}'(a) - g'(a))}_{g''(b) = g''(a) = 0} - \sum_{j=0}^N \int_{x_j}^{x_{j+1}} g'''(x) (\tilde{g}'(x) - g'(x)) dx \\ &= - \sum_{j=0}^N g'''(x_j^+) (h(x_{j+1}) - h(x_j)) = 0, \text{ since } h(x_i) = 0 \text{ at each knot} \end{aligned}$$

Then, we evaluate $\int_a^b \tilde{g}''(t)^2 dt - \int_a^b g''(t)^2 dt = \int_a^b (\tilde{g}''(t) - g''(t))^2 dt - 2 \int_a^b g''(t) (\tilde{g}''(t) - g''(t)) dt$

$$\begin{aligned} &= \int_a^b (\tilde{g}''(t) - g''(t))^2 dt + 2 \int_a^b g''(t) h''(t) dt \\ &= \int_a^b (\tilde{g}''(t) - g''(t))^2 dt \geq 0, \text{ and we conclude that } \int_a^b \tilde{g}''(t)^2 dt \geq \int_a^b g''(t)^2 dt, \end{aligned}$$

and if equal, then $h(t) = 0$ over the interval $[a, b]$, implying $\tilde{g} = g$.

Returning to the penalized least square problem, we see that there are two terms:

$$\min_f \left[\sum_{i=1}^N (y_i - f(x_i))^2 + \lambda \int_a^b f''(t)^2 dt \right]$$

The first term is the data fidelity term, and the second is the penalizing function for curvature. The minimizer is then the function f that has the least penalty for curvature. If we had a set of functions that perfectly interpolate the data, the Natural Cubic spline would be in this set. Since we proved that the natural cubic spline has the minimum curvature in that set, the minimizer is then the Natural Cubic spline with knots @ x_i .

3)

For the question, I implemented K-Nearest Neighbors and Polynomial Regression using the SciKit Learn (SKLearn) library in python to try to minimize the noise in the two functions given in the question. First, I implemented functions that represented the functions provided in the problem for 'bumps' and for 'doppler'. Then, I followed the same equation given by the problem to generate 2048 linearly spaced points between 0 and 1 and executed the function over those values to generate the noise-free function plots. The root loss function was then defined by using the equation given in the problem statement as well. From the graphs provided by Donoho et al., the sigma value of the noise was provided to be 1, and the signal scaling value was estimated based on the graphs provided. The scaling factor for the bumps function was estimated to be 10, and the scaling factor for the doppler function was estimated to be 25. For each replication, Gaussian noise was generated by taking 2048 random values from the standard normal ($z \sim N(0,1)$). Polynomial regression and KNN regressions were fit per the documentation provided by the SKLearn library, and both the root-mean-squared error and the loss were captured for each replication performed. After all replications were completed, the mean of the 20 root loss values were taken for each regression type and each function, as well as the average RMSE as evidenced below. Additionally, the graph of the last run is provided below to also identify the fit.

```
Average RMSE of Root Loss for 20 Replications:
Bumps Function:
    Polynomial Regression: 6.4536
    K-Nearest-Neighbors Regression: 0.8587
Doppler Function:
    Polynomial Regression: 6.4274
    K-Nearest-Neighbors Regression: 0.4890

Average Root Mean Squared Error Over 20 Replications:
Bumps Function:
    Polynomial Regression: 6.4536
    K-Nearest-Neighbors Regression: 0.8587
Doppler Function:
    Polynomial Regression: 6.4274
    K-Nearest-Neighbors Regression: 0.4890
```

Figure 1: Output from the provided python function. The Average RMSE of the Root Loss and Average Root Mean Squared Error for 20 Replications for both types of regression for both functions provided in the homework question. .

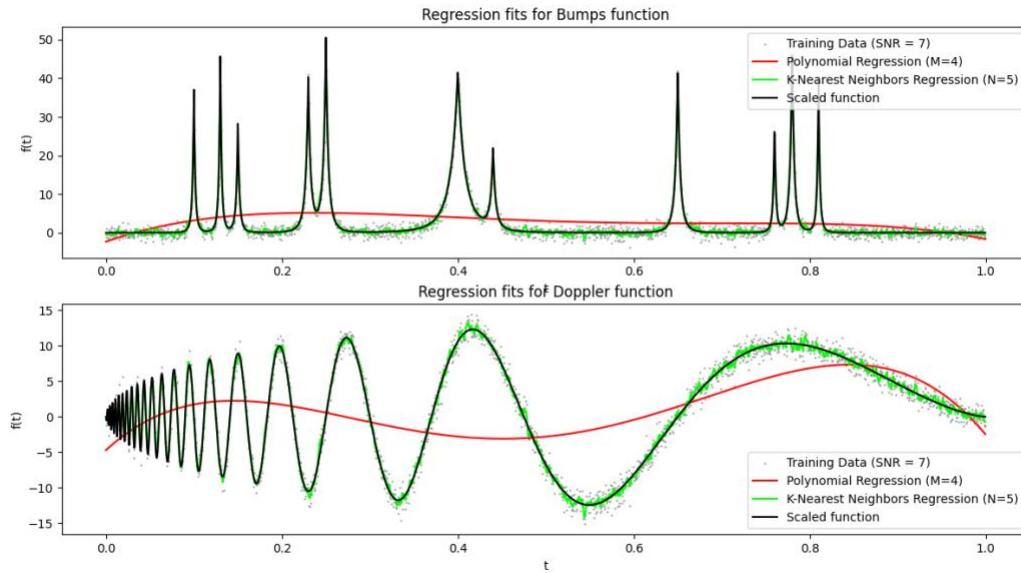


Figure 2: Graphs indicating the fit for both functions. The green line is the KNN regression, and the red line is polynomial regression. The gray dots are the noisy values that the two regression functions were trained and fit against.

To supplement, we see that the polynomial regression does an extremely poor job of estimating both functions, even without taking the noise into account. Neither function was able to be well defined using a polynomial of degree 4, which is why polynomial regression was beat by the wavelet methods in both cases by significant margins. To summarize, the polynomial regression was not the correct regression method to use in this scenario, and thus, the loss value and RMSE values were large.

Regarding the K-Nearest-Neighbors regression with $k=5$, we see that it does a much better job estimating the original function than the polynomial regression method used previously. We see that the loss/RMSE values for the KNN regression were 0.8587 and 0.4890 for the Bumps function and the Doppler function, respectively. For the Bumps function, we see that the RMSE is beaten by all the wavelet methods except for VisuShrink using S8. For the Doppler function, we see that the RMSE of the KNN regression is better than some and worse than some, roughly somewhere in the middle when it comes to efficiency. I believe that wavelet shrinkage was superior to KNN regression in these cases because of the logic behind the algorithms. KNN takes the nearest k neighbors to a point and analyzes the approximate value of what the point should 'look' like, which does not account for noise, meaning that noise would affect the regression function. This is in contrast to wavelet transforms, which splits the data into approximations (lowpass) and detail (highpass) portions, minimizes the highpass terms, and takes the inverse to identify the de-noised signal. In essence, the reason why wavelet shrinkage was better than KNN is because the wavelet shrinkage minimized the effects of noise on the regression, while KNN regression took noise into the regression pattern, and did not account for it.