

Q2

Part 1

(1) Since $m = 0$, which means the f term in the second term should be $f(x)$; $\lambda = \infty$, hence, when we are minimizing \hat{f}_1 , the second term (which penalizes the curvature of the function) should be enforced to be 0, that is $f(x) = 0$. It has zero polynomial (or undefined degree of polynomial) since the function $f(x)$ is actually a zero constant.

(2) $m = 1$ gives us the first derivative $f'(x)$. Similar to (1), $\lambda = \infty$, which means the second term which penalizes curvature of the function should be enforced to zero, that is, $f'(x) = 0$. $f(x)$ is a constant, hence the degree of polynomial is 0. \hat{f}_1 in this case is basically a least square regression.

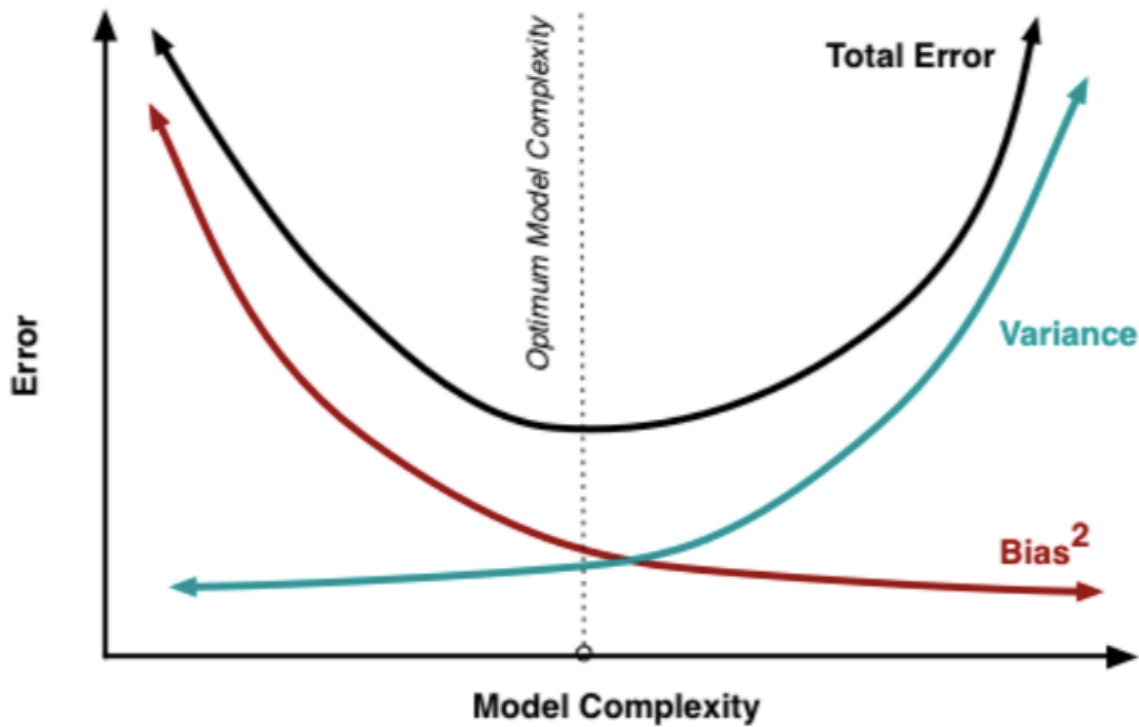
(3) $m = 2$ gives us the second derivative $f''(x)$. $\lambda = \infty$ means when we are minimizing \hat{f}_1 , the second term (which penalizes the curvature of the function) should be enforced to be 0. In this case, $f''(x) = 0$. $f(x)$ is a polynomial of degree 1, takes the form of $f(x) = ax + b$, where a, b are constant. \hat{f}_1 is a linear least square regression in this case.

(4) $m=3$ gives us the third derivative $f'''(x)$. $\lambda = 0$ which basically drop the second term which penalizes the curvature of the function. \hat{f}_1 only tries to minimize the first term which measures the closeness of the model to the data. $f(x)$ will interpolate all data points which leads to overfitting (a low bias but very high variance). $f(x)$ can have an arbitrarily large polynomial degree.

Part 2

(1) When λ goes to infinity, both $f^{(m)}(x)$ and $f^{(m+k)}(x)$ will be enforced to 0. Since k is a positive integer, $f^{(m+k)}(x)$ is a higher order degree derivative, which means the corresponding $f(x)$ function has more degree of polynomial. This leads to \hat{f}_2 to have a higher variance, lower bias. Hence \hat{f}_2 has smaller training RSS.

(2) When λ goes to infinity, the argument is similar to (1), \hat{f}_1 has a lower degree polynomial, hence \hat{f}_1 has higher bias, lower variance. \hat{f}_2 has higher variance, lower bias. But since we are comparing test RSS, according to the bias-variance tradeoff (the plot taken from internet shown below), the test RSS depends on which model has a smaller $\text{bias}^2 + \text{variance}$, which we cannot compare without sufficient conditions here.



(the x axis should be the degree of polynomial in this case)

(3) When $\lambda = 0$, both f_1 and f_2 interpolate data points as much as they can without the constraints of the curvature. \hat{f}_1 and \hat{f}_2 will both overfit (a very low bias but very high variance). \hat{f}_1 and \hat{f}_2 will have the same training RSS, the same test RSS, since the objective functions are basically the same without the second penalization term.

In []: