

Homework 3

Q1

- Suppose that we have a dataset $D = \{(x_i, y_i)\}_{i=1}^n$ of n samples, each having input features $x_i \in \mathbb{R}^d$ and output labels $y_i \in \mathbb{R}$.
- Suppose that there exists a vector $\beta \in \mathbb{R}^d$ (aka. the model) such that $y_i = x_i^\top \beta + \epsilon_i$, for $i \in \{1, \dots, n\}$,

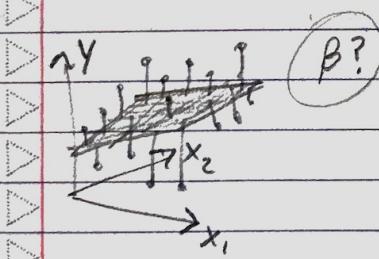
where $\{\epsilon_i\}_{i=1}^n$ are independent and identically distributed Gaussian random variables with mean 0 and variance $\sigma^2 > 0$. Denote by $X = [x_i]_{i=1}^n \in \mathbb{R}^{n \times d}$ the matrix comprising sample features in its rows and $y = [y_i]_{i=1}^n \in \mathbb{R}^n$ the vector of all labels.

$$X \in \mathbb{R}^{n \times d} \quad \begin{matrix} \text{Height} \\ \text{Weight} \\ \text{Age} \\ \text{Blood pressure} \end{matrix} \quad D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$Y \in \mathbb{R}^n \quad y_i = \beta^\top x_i + \epsilon_i, \quad i = 1, \dots, n$$

$$\epsilon_i \sim N(0, \sigma^2) \text{ i.i.d.}$$

$$X \in \mathbb{R}^{n \times d} \quad Y \in \mathbb{R}^n$$



- (1.1) Prove that the maximum likelihood estimator (MLE) equals the least-squares estimator (LSE), i.e.:

$$\hat{\beta}^{ML} = \arg \max_{\beta \in \mathbb{R}^d} P(D|\beta) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \beta^\top x_i)^2$$

- Do we need to know σ^2 to compute $\hat{\beta}^{ML}$?

Say $y_i = f(x_i) + \epsilon_i$ Computing least squares error is like saying you are estimating data with some gaussian noise in i.i.d.

$$(1.1 \text{ continued}) \quad \hat{\beta}^{ML} = \arg \max_{\beta \in \mathbb{R}^n} P(D|\beta) = \arg \min_{\beta \in \mathbb{R}^n} \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

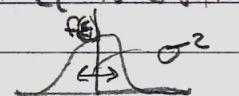
$$P(D|\beta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} [e^{-(y_i - \beta^T x_i)^2 / 2\sigma^2}]$$

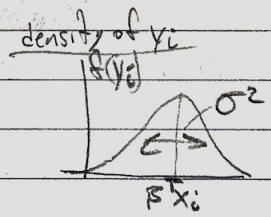
this is a lot to do out so
we can take the log since

$$\arg \max_{\beta \in \mathbb{R}^n} \log P(D|\beta) = \arg \max_{\beta \in \mathbb{R}^n} P(D|\beta)$$

Also, we know this to be the equation
for $P(D|\beta)$ because the noise/error

$\{\epsilon_i\}_{i=1}^n$ are iid. gaussian random variables
so y_i also is of the same form. The
probability of the dataset is given by the
product of a bunch of gaussians.

if $\beta^T x_i = f(x_i)$
we have $y_i = f(x_i) + \epsilon_i$
so if ϵ_i is distributed
like: 



Then y_i will have the
same shape just
shifted by $\beta^T x_i$.
i.e. $y_i \sim N(\beta^T x_i, \sigma^2)$

Now we know why we use the gaussian distribution for the equation for $P(D|\beta)$.

We can take the log as said above but similarly we can minimize the
negative log, $\arg \max_{\beta} P(D|\beta) = \arg \min_{\beta} -\log P(D|\beta)$

$$-\log P(D|\beta) = -\log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-(y_i - \beta^T x_i)^2 / 2\sigma^2}$$

$$= -\log \underbrace{\left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n}_{\text{This term becomes a constant we can call } C} \prod_{i=1}^n e^{-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}}$$

$$= -\sum_{i=1}^n \log e^{-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}} - C$$

$$= \sum_{i=1}^n \frac{(y_i - \beta^T x_i)^2}{2\sigma^2} - \left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T x_i)^2 \right) - C$$

When we find the $\arg \min$ it does not matter if we subtract
a constant C or multiplying by a scalar $\frac{1}{2\sigma^2}$

Homework 3

(1.1 cont) so if $\arg \min_{\beta}$ is the same then we can say

$$\arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

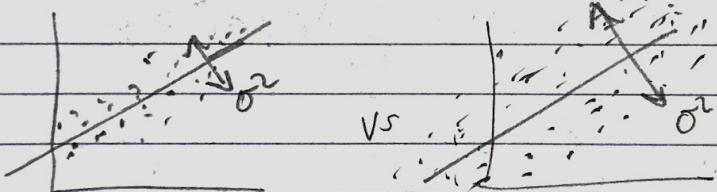
so we arrive at the minimum (or least) of the error
 $(y_i - \beta^T x_i)$ squared.

$$\text{so we finally see, } \hat{\beta} = \arg \max_{\beta \in \mathbb{R}^d} P(\beta | \beta) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

we also see that we do not need σ^2 to compute $\hat{\beta}^{\text{ML}}$!

Since it is a scalar it did not impact our argument.

also intuitively the noise around our "best fit" / least square error line does not impact what the line should be. ex:



(1.2) Show that $E[\hat{\beta}^{\text{ML}}] = \beta$ and $\text{cov}(\hat{\beta}^{\text{ML}}) = \sigma^2 (X^T X)^{-1}$

$$y_i = \beta^T x_i + \epsilon_i, \epsilon_i \text{ are i.i.d., } E[\epsilon_i] = 0, E[\epsilon_i^2] = \sigma^2$$

$i=1, \dots, n$

\hookrightarrow mean of the noise \hookrightarrow variance of the noise

$$y = X\beta + \epsilon$$

$$\hat{\beta}^{\text{ML}} = \arg \max_{\beta \in \mathbb{R}^d} P(\beta | \beta) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

To find the expectation of the ML estimate of β we can

$$\text{also look at the square error loss } f(\beta) = \sum_{i=1}^n (y_i - \beta^T x_i)^2.$$

setting its gradient equal to zero will help us

find the $\arg \min_{\beta}$ and we will use that to find $E[\hat{\beta}^{\text{ML}}]$

(1.2 continued)

$$\nabla f(\beta) = \sum_{i=1}^n \nabla f_i(\beta), \text{ where } f_i(\beta) = (y_i - \beta^T x_i)^2$$

$$\nabla f_i(\beta) = \begin{bmatrix} \frac{\partial f_i(\beta)}{\partial \beta_1} \\ \vdots \\ \frac{\partial f_i(\beta)}{\partial \beta_n} \end{bmatrix} \in \mathbb{R}^d \text{ and } \frac{\partial f_i(\beta)}{\partial \beta_k} = \frac{\partial}{\partial \beta_k} (y_i - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_d x_{id})^2$$

$\underbrace{}$

$$= \frac{\partial}{\partial \beta_k} (y_i - \underbrace{\beta_k x_{ik}}_{\beta^T x_i} - \sum_{j \neq k} \beta_j x_{ij})^2$$

$$= -2(y_i - \beta^T x_i)x_{ik}$$

$$\text{and so, } \nabla f_i(\beta) = -2(y_i - \beta^T x_i)x_i \in \mathbb{R}^d$$

$\underbrace{\beta^T x_i}_{\in \mathbb{R}} \quad \underbrace{x_i \in \mathbb{R}^d}$

$$\begin{aligned} \text{also, } \nabla f(\beta) &= 2 \sum_{i=1}^n (\beta^T x_i - y_i)x_i = 2 \sum_{i=1}^n x_i (x_i^T \beta - y_i) \\ &= 2 \sum_{i=1}^n x_i x_i^T \beta - 2 \sum_{i=1}^n x_i y_i \\ &= 2 X^T X \beta - 2 X^T y \end{aligned}$$

$$X = \begin{bmatrix} x_1 & \dots \\ x_2 & \dots \\ \vdots & \ddots \\ x_n & \dots \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$\arg \min_{\beta} \text{ where } \nabla f(\beta) = 0 \text{ so, } 0 = 2 X^T X \beta - 2 X^T y$$

$$2 X^T y = 2 X^T X \beta$$

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

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y, \quad y = X \beta + \varepsilon \\ \hat{\beta} &= (X^T X)^{-1} X^T (X \beta + \varepsilon) \\ &= (X^T X)^{-1} X^T X \beta + (X^T X)^{-1} X^T \varepsilon \\ &= I \beta + G \varepsilon \end{aligned}$$

$$\hat{\beta} = I \beta + G \varepsilon$$

$$\mathbb{E}[\hat{\beta}] = \beta + \mathbb{E}[G \cdot \varepsilon] = \beta + G \cdot \mathbb{E}[\varepsilon] = \beta + G \cdot \mathbf{0} \quad \Leftrightarrow$$

$$\mathbb{E}[\hat{\beta}] = \beta$$

sum of the expectations is the expectation of the sums

Homework 3

$$(1.2 \text{ cont.}) \text{ cov}(\hat{\beta}^{\text{ML}}) = \sigma^2 (X^T X)^{-1}$$

$$\text{cov}(\hat{\beta}^{\text{ML}}) = E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])^T], \text{ outer product} \Rightarrow \text{matrix } \in \mathbb{R}^{d \times d}$$

$= E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T]$, we have shown $\hat{\beta} = \beta + G\epsilon$ where

$$[\text{so, } \hat{\beta} - \beta = G\epsilon]$$

$$G = (X^T X)^{-1} X^T$$

$$= E[G\epsilon \cdot \epsilon^T (G^T \cdot \epsilon)^T], \text{ where } \epsilon \text{ is the vector containing all of the noise}$$

$$= E[G \epsilon \epsilon^T G^T], \text{ since } G \text{ is a constant then}$$

$$= G E[\epsilon \cdot \epsilon^T] G^T$$

$$\left[E[\epsilon \cdot \epsilon^T] = \begin{bmatrix} E[\epsilon_1^2] & E[\epsilon_1 \epsilon_2] & \dots \\ \vdots & \ddots & \vdots \\ E[\epsilon_n^2] \end{bmatrix} \right] = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \sigma^2 \end{bmatrix}, E[\epsilon_i \epsilon_j] = E[\epsilon_i] E[\epsilon_j] \xrightarrow{\text{independent}} \text{so, } = 0$$

$$= \sigma^2 I$$

$$= G \sigma^2 I G^T$$

$$= \sigma^2 G G^T$$

$$= \sigma^2 (X^T X)^{-1} X^T [(X^T X)^{-1} X^T]^T$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} \cdot I$$

$$\text{cov}(\hat{\beta}^{\text{ML}}) = \sigma^2 (X^T X)^{-1}$$

Note: if the noise is iid. gaussian then

$$\text{our } \hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

unbiased estimate is variable w/ σ^2 .

The bigger σ^2 , the bigger the noise of y s, the bigger the covariance of $\hat{\beta}$. The bigger the noise the more uncertain we are of β . Our estimate becomes worse.

(1.3) Consider now a new sample $(x_0, y_0) \in \mathbb{R}^d \times \mathbb{R}$ outside the dataset, for which: $y_0 = x_0^\top \beta + \varepsilon_0$
 where $\varepsilon_0 \sim N(0, \sigma^2)$, and is independent of all the noise variables $\{\varepsilon_i\}_{i=1}^n$ in D . Show that the expectation prediction error (EPE) of estimate $\hat{y}_0 = x_0^\top \hat{\beta}^{\text{ML}}$ is given by: $E[(y_0 - \hat{y}_0)^2] = \sigma^2 + x_0^\top \text{cov}(\hat{\beta}^{\text{ML}}) x_0$

EPE of $\hat{y}_0 = x_0^\top \hat{\beta}^{\text{ML}}$ comes from our noise ε_0 and we don't know β
 so we estimate it, $\hat{\beta}^{\text{ML}}$. ε_0 is irreducible error
 $\hat{\beta}^{\text{ML}}$ has estimation error.

$$x_0 \in \mathbb{R}^d, \quad y_0 = x_0^\top \beta + \varepsilon_0, \quad E[\varepsilon_0] = 0, \quad E[\varepsilon_0^2] = \sigma^2, \quad \varepsilon_i \text{ are i.i.d.}$$

$$\hat{y}_0 = x_0^\top \hat{\beta}^{\text{ML}}$$

So,

$$E[(y_0 - \hat{y}_0)^2] = E[(y_0 - x_0^\top \hat{\beta}^{\text{ML}} + x_0^\top \hat{\beta}^{\text{ML}} - \hat{y}_0)^2]$$

$$= E[\underbrace{(y_0 - x_0^\top \hat{\beta})^2}_A + \underbrace{(x_0^\top \hat{\beta} - \hat{y}_0)^2}_B + 2(y_0 - x_0^\top \hat{\beta}) \cdot \underbrace{(x_0^\top \hat{\beta} - \hat{y}_0)}_C]$$

$$A = (y_0 - x_0^\top \hat{\beta})^2 = \varepsilon_0^2$$

$$E[A] = E[\varepsilon_0^2] = \sigma^2$$

$$C = (y_0 - x_0^\top \hat{\beta})(x_0^\top \hat{\beta} - \hat{y}_0) = \varepsilon_0(x_0^\top \hat{\beta} - \hat{y}_0)$$

$$= \varepsilon_0(x_0^\top \beta - x_0^\top \hat{\beta}^{\text{ML}}), \quad \hat{\beta}^{\text{ML}}$$
 is random since our data is random

$$E[C] = E[\varepsilon_0] E(x_0^\top \beta - x_0^\top \hat{\beta}^{\text{ML}}) = 0 \cdot 0 = 0$$

$$\text{for anti-}\vec{y} \text{ vectors}$$

$$x^\top \vec{y} = \vec{y}^\top x$$

$$B = (x_0^\top \hat{\beta} - \hat{y}_0)^2$$

$$E[B] = E[(x_0^\top \hat{\beta} - x_0^\top \hat{\beta}^{\text{ML}})(\hat{\beta}^\top x_0 - \hat{\beta}^{\text{ML}}^\top x_0)]$$

$$= E[x_0^\top (\beta - \hat{\beta}^{\text{ML}})(\hat{\beta}^\top - \hat{\beta}^{\text{ML}})^\top x_0]$$

$$= x_0^\top E[(\beta - \hat{\beta}^{\text{ML}})(\hat{\beta}^{\text{ML}} - \beta)^\top] x_0$$

$$= x_0^\top E[\text{cov}(\hat{\beta}^{\text{ML}})^\top] x_0$$

$$= x_0^\top \text{cov}(\hat{\beta}^{\text{ML}}) x_0$$

note: Estimator and covariance is undefined if $\text{rank}(X) < d$

$$E[(y_0 - \hat{y}_0)^2] = E[A+B+C] = \sigma^2 + x_0^\top \text{cov}(\hat{\beta}^{\text{ML}}) x_0$$

off by, error noise + estimation error

Homework 3

Q2

- ▷ Consider the exact same setting as Q1, but suppose now that we further assume that $\beta \sim N(0, \frac{\sigma^2}{\lambda} X^T I)$, i.e. that the model β is sampled from a Gaussian prior centered at zero, with a covariance that depends on $\lambda > 0$.
- ▷ 2.1) What information does the prior capture / what does it tell us about our prior belief on β as $\lambda \rightarrow \infty$?
- ▷ As $\lambda \rightarrow \infty$, the variance ($\frac{\sigma^2}{\lambda}$) becomes very small which indicates a high confidence in our mean (0). It tells us our prior belief on β is low; that there is low correlation between X and y , (our features and classes).
- ▷ 2.2) What information does the prior capture / what does it tell us about our prior belief on β as $\lambda \rightarrow 0$? What did we call such a prior in previous lectures?
- ▷ As $\lambda \rightarrow 0$ the variance ($\frac{\sigma^2}{\lambda}$) becomes very large and approaches the uninformative prior. $I +$ would be uninformative if all of our coefficients in β were 1. But it indicates a high correlation between all of our features and the classifications.
- ▷ (2.3) Show that the maximum a posteriori estimate corresponds to ridge regression i.e.:

$$\begin{aligned}\hat{\beta}^{MAP} &\equiv \arg \max_{\beta \in \mathbb{R}^d} P(\beta | D) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|_2^2 \\ \arg \max_{\beta} P(\beta | D) &= \arg \max_{\beta \in \mathbb{R}^d} P(D | \beta) P(\beta) \\ &= \arg \min_{\beta \in \mathbb{R}^d} -\log P(D | \beta) - \log P(\beta) \\ &= \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|_2^2\end{aligned}$$

(2.4) Where does $\hat{\beta}^{\text{MAP}}$ converge to as $\lambda \rightarrow \infty$?

For a gaussian prior centred around zero, the variance would shrink very small so you would have a β super concentrated around zero.

(2.5) Where does $\hat{\beta}^{\text{MAP}}$ converge to as $\lambda \rightarrow 0$? Contrast this to what happens to the prior: do we see a phenomenon that we have also observed in other parameter estimation tasks?

$\hat{\beta}^{\text{MAP}}$ converges to the uninformative prior and then we also see that $\hat{\beta}^{\text{MAP}} = \hat{\beta}^{\text{ML}}$.

This is also apparent when we did our proof before,

$$\hat{\beta}^{\text{ridge}} = \hat{\beta}^{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|_2^2$$

as $\lambda \rightarrow 0$ we have

$$\arg \min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + 0$$

which we recognize from before as the maximum likelihood estimate; $\hat{\beta}^{\text{ML}}$

So we can see ridge regression as an interpolation between Map and MLE and λ as an interpolation between an uninformative prior and a prior that implements Occam's razor (one that implements the KISS principle \rightarrow prefer it simple)

Homework 3

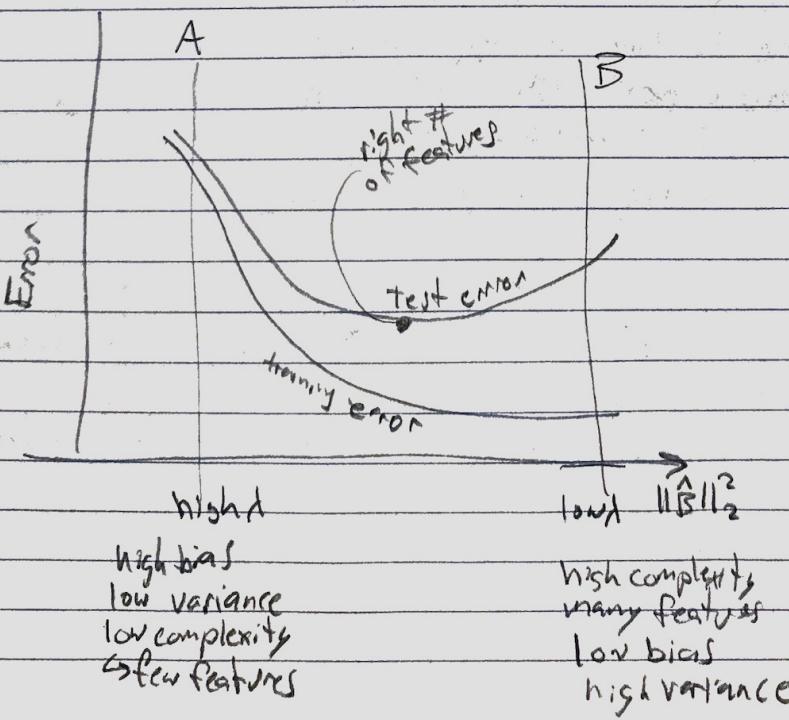
Q3

(3.1) For what value of λ is the training error minimized and why?

$\lambda = 0$, Training error is minimized at $\lambda = 0$ since it corresponds to the most complex model where the model is fitted to the training data set's noise. Also the error is $(y_i - \beta^T x_i)^2 + \lambda \|\beta\|_2^2$, so $\lambda = 0$ minimizes the error in that equation since $\lambda \|\beta\|_2^2$ would become zero.

(3.2) How does the norm of the estimated model $\|\hat{\beta}\|_2$ change as we increase λ and why?

As we increase λ , the norm decreases as seen on the graph. More features are zero so the model is less complex and the norm of the estimated model is smaller.



3.3) Between points A and B, which one corresponds to a model trained with a higher value of λ ? A

3.4) Between A, B, which are we learning more complex model, more parameters away from 0? B

3.5) Between A, B which has higher test/estimation error? B

3.6) higher model error / bias? A

3.7) Which is overfitting the training set? B

3.8) If we increase the number of samples in the training set, would the test error curve change and how?

Model would be better so test error minimum would decrease. More samples will decrease the error more near B than A. Near A the model is simple, has high bias and more data will not fix that but near B, the model error is low and the model will be better when fed by more data. More data will inform the model when trained, which coefficients are zero such as Zipcode when considering blood pressure, so the test error will decrease.

3.9) If we decrease the noise variance σ^2 , would the test error curve change and how? If you decrease the noise the model will be trained more accurately so the test error will decrease, and be more closely fitted to the training error.

Homework 3

Q4

- ▷ (4.1) Provide the mathematical formula for function fun in generate_data.py.
- ▷ What coordinates of vector $x \in \mathbb{R}^5$ does its output depend on?
- ▷ Is it a linear function of $x \in \mathbb{R}^5$?

$$\text{fun}(x) = 1.3 + 2(x_0) - 1.1(x_1) + 0.7(x_2) + 1.2(x_3) + 0.4(x_0^2) - 1.5(x_1)(x_3) - 0.7(x_4^2)$$

The output of $\text{fun}(x)$ depends on x_0, x_1, x_2, x_3, x_4

It is not a linear function since it has terms x_0^2, x_4^2

- ▷ (4.2) Run code on generated data with $N=1000, d=5, \sigma=0.01, fr=10\% \rightarrow 100\%$ 10% increments. Plot train+test RMSE vs # of train samples given to the model.

Does increasing the number of samples significantly reduce either errors?

Do you think that increasing the number of samples beyond 1000

would help in significantly lowering either error? If not, why?

Increasing the number of samples does not significantly reduce either error, so no increasing the number of samples beyond 1000 would not help to significantly lower either error. The model is linear, but the generated data is from a polynomial, so the model error is going to always outweigh the noise of the data, so more data will not fix the model error.

* see plot

(4.3) Run code on a dataset generated by generate_data.py with $N=1000$, $d=40$, $\sigma=0.01$. Note, increasing input vector adds variables that do not affect labels; often called nuisance or redundant variables. Plot train and test RMSE vs training sample size. Do you see any differences from Q4.2? Explain what you think is happening to cause this behavior. Do you think that increasing the number of samples beyond 1000 would significantly reduce the test error?

The main difference from Q4.2 plot is the test error increased, also the plots are less scattered and settle quickly to an error/RMSE. Since the model is trained on more categories $x_5 - x_{39}$ but the generated data has not changed the function or variables the function depends on ($x_0 - x_4$) then the model is taking into account random data that does not actually matter. So this introduces more noise. $\beta_5 - \beta_{39}$ does decrease with more samples but the model is still trying to fit linear to a polynomial so the bulk of the error is still there and will not be fixed with more samples, so beyond 1000 would not significantly reduce test error.

* 4.3 see plot

Homework 3

4.4) see code for def lift(x_initial)

which expands or 'lifts' x to x' (vectors)

by appending x_i and $x_i \cdot x_j$ for all $i \geq j$ to x'

4.5) Show that function: fun() can be written as a linear function of such a lifted vector x' , i.e. a function of the form:

$$f(x') = \beta^T x' + c$$

for appropriately defined vector $\beta \in \mathbb{R}^{d'}$ and constant $c \in \mathbb{R}$, and provide β and c .

$$\text{fun}(x) = 1.3 + 2(x_0) - 1.1(x_1) + 0.7(x_2) + 1.2(x_3) + 0.4(x_4) - 1.5(x_1)(x_2) - 0.7(x_2)^2$$

$$\begin{aligned} x' = & [x_0, x_1, x_2, x_3, x_4, \\ & x_0^2, x_1 x_0, x_1^2, x_2 x_0, x_2 x_1, \\ & x_2^2, x_3 x_0, x_3 x_1, x_3 x_2, x_3^2, \\ & x_4 x_0, x_4 x_1, x_4 x_2, x_4 x_3, x_4^2] \end{aligned}$$

$$f(x') = 2(x'_0) - 1.1(x'_1) + 0.7(x'_2) + 1.2(x'_3) + 0.4(x'_4) - 1.5(x'_1)(x'_2) - 0.7(x'_2)^2 + 1.3$$

$$f(x') = \beta^T x' + c \quad \text{where } c = 1.3 \text{ and,}$$

$$\begin{aligned} \beta^T = & [2, -1.1, 0.7, 1.2, 0, \\ & 0.4, 0, 0, 0, 0, \\ & 0, 0, -1.5, 0, 0, \\ & 0, 0, 0, 0, -0.7] \end{aligned}$$

4.6) My code for lift(x-initial) is actually a step ahead and does what liftDataset should do.

- ▷ (4.7) Lift the X 's in the data set by adding function to Linear Regression SweepN w/ $N=1000$, $d=5$, $\sigma=0.01$.
- ▷ How does this plot differ from previous plots you generated?
- ▷ Why? Obtain the parameters (coefficients and intercept) of your finally trained model. Did you manage to learn function fun?

This plot converged to $\text{RMSE} \approx 0.01$ with test RMSE slightly higher. this is a much better result than before since with the lifted parameters we can now have an accurate model and decrease the model error especially with larger data sets.

The script was able to learn the function fun with the coefficients only ± 0.001 off of the actual values.

$$\text{cx } \beta_0 = 1.99994 \approx 2.0$$

▷ seeplot

- ▷ (4.8) Repeat but with $d=40$. How does the training and test error change as you increase the number of training samples? Obtain the parameters (coefficients and intercept) of your finally trained model. Did you manage to learn function fun? Do you think that increasing the number of samples beyond 1000 would improve the test error / the proximity of your learned model to the true fun fun?

The training error increases with increasing training samples but, The training error is very low compared to the test error. The test error decreases as you increase the number of training samples. The script got close to learning the function fun, off by about ± 0.3 for the coefficients, with the zero coefficient also off, not quite zero. Yes, increasing the number of samples beyond 1000 would improve the test error as seen in the trend already, and made possible by better matching the true function fun, so the proximity of the learned model would improve.

▷ see plots

Bonus Question

LassoCV.py?

1. Read through the code in file LassoCV.py and explain what it does.

It loads the X, y data the same, and splits it into test/train the same

assigns alpha=0.1, It also imports Lasso, from sklearn model

creates a lasso model w/ alpha, where alpha acts as our λ on L₁ norm
assigns cv to a KFold with 5 splits, and shuffle as true

creates a score variable using cross_val_score function with
variables / inputs, model, X_train, y_train, and sets cv
to cv, and scoring to the negative RMSE.

Prints, cross-validation RMSE for $\alpha = \text{alpha}$ value and then
the negative of the average of the scores as well
as the standard deviation of the scores.

Tells you when its done fitting the linear model.
and fits i) the X training and y training data
to the lasso model

Finally computes the RMSE for the train and test classifications
with inputs of y and the models predictions of ys to
see how far off the model is.

and Finally prints the training and testing RSME

Also cross_val_score splits the dataset X and y into cross validation
folds $cv=n$ folds, the estimator is trained and the score
is recorded. After doing it cv times, cross_val_score returns
an array of scores one for each fold. Then we can calculate
the mean and std to get an overall estimate of the models performance.

(Bonus 3) Train RMSE = 0.009983

Test RMSE = 0.0108

best $\alpha = 0.000977$

$\beta_0 = 1.99912 \quad \beta_1 = 0.39983$ all close to fun

Yes! it was better since k-fold and sweeping was able to determine which variables should be zero.
(coeffr)

see plot

git repo @ barker-ch

IML

HW3