

Homework 8 MATH531

Doug Nychka

2025-03-09

Getting started

Loading the AudiA4 data and creating the X matrix for “broken” line regression.

```
library( fields) # load fields package
```

```
## Loading required package: spam
```

```
## Spam version 2.11-1 (2025-01-20) is loaded.  
## Type 'help( Spam)' or 'demo( spam)' for a short introduction  
## and overview of this package.  
## Help for individual functions is also obtained by adding the  
## suffix '.spam' to the function name, e.g. 'help( chol.spam)'.
```

```
##
```

```
## Attaching package: 'spam'
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      backsolve, forwardsolve
```

```
## Loading required package: viridisLite
```

```
##
```

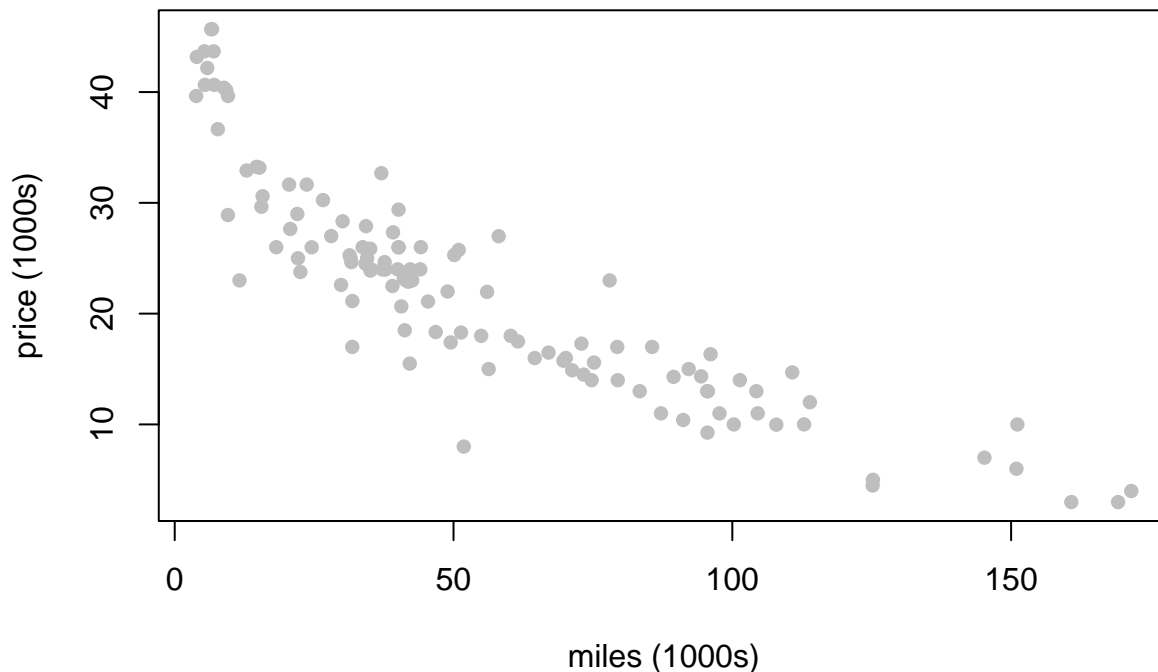
```
## Try help(fields) to get started.
```

```
setwd("~/Dropbox/Home/Teaching/MATH531/MATH-531/MATH531S2024/Homework")  
load("AudiA4.rda" )  
head( AudiA4)
```

```
##      year price mileage distance  
## 58  2020 39649    3848         29  
## 145 2020 43175    3962          7  
## 10  2020 43675    5316          7  
## 52  2020 40649    5417         29  
## 143 2020 42175    5846          7  
## 9   2020 45675    6539          7
```

```
# change units so the numbers are simpler
price<- AudiA4$price/1000 # in thousands of dollars
mileage<- AudiA4$mileage/1000 # thousands of miles

# the data
plot(mileage, price, col="grey", pch=16,
     xlab= "miles (1000s)", ylab="price (1000s)")
```



Problem 1

Here is a standard OLS fit to these data. (You might want to review the handy “I” syntax in an lm formula.)

```
OLSFit<- lm( price~ mileage + I(mileage^2))
summary( OLSFit)

##
## Call:
## lm(formula = price ~ mileage + I(mileage^2))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.7513  -1.9017  -0.3323   2.2942   9.1094
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  39.4106402  0.9287195  42.435  < 2e-16 ***
## mileage      -0.4350211  0.0299689 -14.516  < 2e-16 ***
## I(mileage^2)  0.0014482  0.0001927   7.514  1.3e-11 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.968 on 116 degrees of freedom
## Multiple R-squared:  0.8422, Adjusted R-squared:  0.8395
## F-statistic: 309.6 on 2 and 116 DF,  p-value: < 2.2e-16
```

- (a) Use a qqplot to assess if the standardized residuals follow a normal distribution, $N(0,1)$ To do this from base R it is

```
n<- length(OLSPfit$residuals)
p<- length(OLSPfit$coefficients)
sigma2Hat<- sum(OLSPfit$residuals^2)/ ( n-p)
theoretical<- qnorm( ((1:n) -.5)/n )

plot( theoretical, sort(OLSPfit$residuals/sqrt(sigma2Hat))
)
abline( 0,1, col="magenta")
```

As part of your assessment add 90% bounds at each theoretical quantile for what you would expect if the data was drawn from iid $N(0,1)$.

Hint: These 90% intervals will involve generating Monte Carlo samples and then finding 5% and 95% quantiles. Also generate simultaneous intervals using the Bonferroni adjustment: $.05/n$ and $1 - .05/n$.

Problem 2

Create a set of basis functions as bumps with the form:

$$H(d) = (1 + d)e^{\ell - d}$$

and the i^{th} basis function being

$$\phi_i(u) = H((u - v_i)/\alpha)$$

where $\{v_i\}$ are a grid of values and α a scaling parameter. I use $\alpha = 10$ below, don't change that.

Use the code below to create a matrix where all the basis functions are evaluated at all the mileages.

```
mGrid1<- seq( 0, 175, length.out=50)
bigD<- rdist( mileage, mGrid1)/10
# 10 is a scaling for the width of the bumps
Phi<- (1+ bigD)* exp( -bigD)
```

- 2(a) There are 50 basis functions in this example. Plot the first, 20th and the 45th basis functions over the range of the mileage. Put these on the same figure for comparison.

Hint: These are the columns of Φ .

- 2(b) Fit an OLS model according to

```
FitBasis<- lm( price ~ Phi)
```

Make a scatterplot of the data and add this fitted curve to it. Plot the estimated curve at the finer grid of points `mGrid2<- seq(0, 175, length.out= 250)`

2(c) Now consider a ridge regression estimator:

$$\hat{\beta} = (\Phi^T \Phi + \alpha I)^{-1} \Phi^T y$$

where $\alpha \geq 0$ and y in this case is the price. Predicted values are

$$\hat{y} = \Phi \hat{\beta}$$

Below is a handy function to do this. Note that it is hardwired for this data set and basis.

```
mySmoother<-function(alpha){
  smootherCoef<- solve(t(Phi)%*%Phi + alpha* diag(1,50))%*%t(Phi)%*%price
  ridgeFit<- Phi%*%smootherCoef
  # note smoother Matrix is
  # S<- Phi%*%solve(t(Phi)%*%Phi + alpha* diag(1,50))%*%t(Phi)
  return( ridgeFit)
}
```

and as a code example

```
fitTest<- mySmoother( .001)
plot( mileage,price)
lines( mileage, fitTest)
```

Vary α and choose a value that subjectively looks like a good fit to these data. Add this curve to your figure in 2(b). Also report the “effective number of parameters” in your choice as the trace of the matrix `smootherMatrix`. (You will have to adapt/hack the code for the `mySmoother` function to get this.)

Hint: It is useful to vary α *equally spaced on a log scale* to get different amounts of smoothing. I used an α range of $1e-6$ to $1e2$.

EXTRA CREDIT: Explain how to modify this estimator to include a constant and linear function where as $\alpha \rightarrow \infty$ the ridge estimate is just the OLS estimate of a line. (This should help with getting a better fit at the ends.)

Problem 3

This problem compares the OLS fit quadratic function to the Bayesian version. To make the uncertainty of the parameters more comparable in size use the X matrix:

```
X<- cbind( 1, mileage/10, (mileage^2)/1000 )
```

Because this is a linear transformation you should get the same predicted values and the inferences will be the same. Of course the coefficients are different.

```
OLSFit2<- lm( price ~ X -1)
summary(OLSFit2)$coefficients
```

```
##      Estimate Std. Error    t value    Pr(>|t|)
## X1 39.410640  0.9287195  42.435462 1.704583e-72
## X2 -4.350211  0.2996893 -14.515738 7.575508e-28
## X3  1.448150  0.1927143   7.514493 1.298652e-11
```

Below is an excerpt from the wikipedia page on the Bayesian linear model that details the

2/18/25, 12:29 PM

Bayesian linear regression - Wikipedia

Therefore, the posterior distribution can be parametrized as follows.

$$\rho(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}, \mathbf{X}) \propto \rho(\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}, \mathbf{X}) \rho(\sigma^2 \mid \mathbf{y}, \mathbf{X}),$$

where the two factors correspond to the densities of $\mathcal{N}(\boldsymbol{\mu}_n, \sigma^2 \boldsymbol{\Lambda}_n^{-1})$ and Inv-Gamma (a_n, b_n) distributions, with the parameters of these given by

$$\boldsymbol{\Lambda}_n = (\mathbf{X}^\top \mathbf{X} + \boldsymbol{\Lambda}_0), \quad \boldsymbol{\mu}_n = (\boldsymbol{\Lambda}_n)^{-1} (\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} + \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0),$$

$$a_n = a_0 + \frac{n}{2}, \quad b_n = b_0 + \frac{1}{2} (\mathbf{y}^\top \mathbf{y} + \boldsymbol{\mu}_0^\top \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0 - \boldsymbol{\mu}_n^\top \boldsymbol{\Lambda}_n \boldsymbol{\mu}_n).$$

which illustrates Bayesian inference being a compromise between the information contained in the prior and the information contained in the data.

posterior distribution.

For the priors applied to the quadratic regression model, set $\mu_0 = 0$ $\Lambda_0 = .01$ and for the Inverse gamma prior on σ^2 use $a_0 = 1/2$ and $b_0 = (1/2) * 20$. These will give a prior distribution around 20 with a large spread.

Now sample from the joint posterior 10000 times. That is for 10000 repetitions first sample σ^2 from its posterior IG distribution (IG(a_n, b_n)) and then sample $\boldsymbol{\beta}$ from a multivariate normal conditional on the value sampled for σ^2 Find the mean and standard deviations for the three regression parameters from these 10000 samples and compare them to the OLS estimates and standard errors obtained from the OLS fit above.