

Bayesian Learning (732A73) Lab 2

Hoda Fakhrzadehjahromy (hoda7840), Ravinder Alta (ravat601)

Linear and Polynomial Regression

In [1]: library(formatR)

1. Linear and polynomial regression

The dataset TempLinking.txt contains daily average temperatures (in degree Celcius) in Malmsslätt, Linköping over the course of the year 2018. The response variable is temp and the covariate is the number of days since the beginning of the year

$$time = \frac{\text{the number of days since beginning of year}}{365}$$

A Bayesian analysis of the following quadratic regression model is to be performed:

$$temp = \beta_0 + \beta_1 \cdot time + \beta_2 \cdot time^2 + \varepsilon, \varepsilon \stackrel{iid}{\sim} N(0, \sigma^2)$$

(a)

(a) Use the conjugate prior for the linear regression model. The prior hyper-parameters μ_0, ν_0, ν_0 and σ_0^2 shall be set to sensible values. Start with $\mu_0 = (-10, 100, -100)^T$, $\Omega_0 = 0.01 \cdot I_3$, $\nu_0 = 4$ and $\sigma_0^2 = 1$. Check if this prior agrees with your prior opinions by simulating draws from the joint prior of all parameters and for every draw compute the regression curve. This gives a collection of regression curves; one for each draw from the prior. Does the collection of curves look reasonable? If not, change the prior hyperparameters until the collection of prior regression curves agrees with your prior beliefs about the regression curve. (Hint: R package mvtmnorm can be used and your Inv- χ^2 simulator from Lab 1.)

In [2]: TempLink = read.table("~/TempLinking.txt", header = TRUE)
attach(TempLink)

In [3]: dim(TempLink)

365 2
%%

In [4]: head(TempLink)

```
A data.frame: 6 x 2
  time      temp
  <dbl>      <dbl>
1  0.002740   2.0083
2  0.005479   2.8667
3  0.008219   2.0750
4  0.010959   2.0708
5  0.013699   0.5583
6  0.016438  -3.5208
```

```
# setting the initial values
mu.0 = c(-10,100,10)
omega.0 = 0.01*diag(3)
nu.0 = 4
sigma2.0 = 1
```

```
tau2<- function(data,mu,n){
  sum((data)-(mu)^2)/n
}
```

```
# Random generation from a scaled inverse chisquare
rinvschsq <- function(draws, n, tau) {
  ch1_square <- rchisq(draws, n)
  return( tau*(n-1)/ch1_square )
}
# Density of a scaled inverse chisquare
dinvschsq <- function(data, df, tau) {
  return( tau*2*df/2*(df/2)/gamma(df/2) * exp(-df*tau/2*(2+data)) / data^(1+df/2) )
}
```

In [7]: lmTemp = lm(temp ~ time + I(time^2), data = TempLink)

```
summary(lmTemp)
```

Call:
lm(formula = temp ~ time + I(time^2), data = TempLink)

Residuals:

Min	1Q	Median	3Q	Max
-14.5949	-3.2275	0.0759	3.5015	14.2577

Coefficients:

Estimate	Std. Error	t value	Pr(> t)			
(Intercept)	-11.956	0.820	-14.58	<2e-16 ***		
time	103.584	3.776	27.43	<2e-16 ***		
I(time^2)	-95.418	3.647	-26.16	<2e-16 ***		
---	---	---	---	---		
Signif.	0 ***	0.001 ***	0.01 ***	0.05 **	0.1 *	1

Residual standard error: 5.193 on 362 degrees of freedom
Multiple R-squared: 0.6759, Adjusted R-squared: 0.6741
F-statistic: 377.5 on 2 and 362 DF, p-value: < 2.2e-16

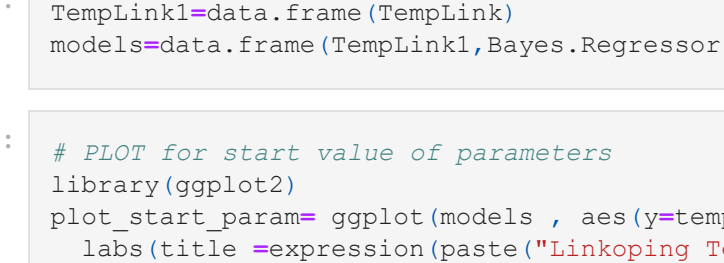
In [9]: sqrt(26.7)

5.16720427310553

In [10]: sum(lmTemp\$residuals**2)/length(lmTemp\$residuals)

26.747622942917

In [11]: plot(y=temp, x=time, col='deeppink', pch=19, lwd=3)
lines(y=lmTemp\$fitted.values, x=time, col='blue', lwd=3)



slide.jpg

In [12]: library(mvtmnorm)

we first simulate σ^2 from its marginal prior Inv- χ^2 and then simulate beta from its prior conditional distribution $N(\nu_0/\Omega_0, \Omega_0/2) | \Omega_0, \Omega_0^{-1} \sim I$

In [13]: sigma2.prior <- function(){
 rinvschsq(draws = 1, n = nu.0, tau = sigma2.0)
}

In [14]: Beta.prior2 <- function(sigma2){
 rvmnorm(mean = mu.0, n = 1, sigma = sigma2*solve(omega.0))
}

In [15]: # create empty structure for sigma, Beta and error
Ndraws = 200
ErrorTerm = numeric(Ndraws)
sigma2 = numeric(Ndraws)
BetaList = matrix(Ndraws, 3)
colnames(BetaList) = c("B0", "B1", "B2")

In [16]: for(i in 1:Ndraws){
 sigma2[i] = sigma2.prior()
 BetaList[i,1] = Beta.prior(sigma2[i])[1]
 BetaList[i,2] = Beta.prior(sigma2[i])[2]
 BetaList[i,3] = Beta.prior(sigma2[i])[3]
 ErrorTerm[i] = rnorm(1, mean = 0, sd = sqrt(sigma2))
}

In [17]: Bayes.Regressor = matrix(, length(time), Ndraws)

In [18]: for(i in 1:Ndraws){
 Bayes.Regressor[i,1] = BetaList[1,i] + BetaList[2,i]*time + BetaList[3,i]*(time^2) +
 colnames(Bayes.Regressor) = paste0("model", 1:Ndraws)
}

In [19]: head(data.frame(Bayes.Regressor), 1)

```
A data.frame: 1 x 10
  model1 model2 model3 model4 model5 model6 model7 model8 model9 model10 ...
  <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl> ...
1 -10.0377 -16.09405 8.828533 -10.42248 -9.4578 -8.270855 3.317529 -5.906997 -21.34576 -5.579045 ...
```

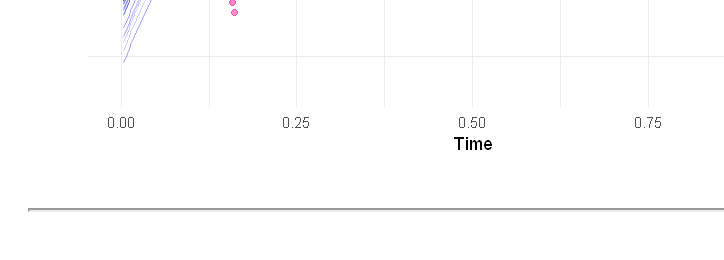
In [20]: TempLink1 = data.frame(TempLink)
models = data.frame(TempLink1, Bayes.Regressor)

In [21]: # Plot for start value of parameters
library(ggplot2)
plot_start_param = ggplot(models, aes(y=temp, x = time)) +
 labs(title = expression(paste("Linkoping Temperature", " ", mu[0], " ", sigma[0], " ",
 Omega[0], " ", nu[0])), x = "Time", y = "Temperature") + theme_minimal()

for(i in names(models2)[-c(1,2)]){
 plot_start_param = plot_start_param +
 geom_line(aes_string(y = i), color = "blue", alpha = 0.2)
}

plot_start_param = plot_start_param +
 geom_point(aes(y = temp), alpha = 0.5, color = "deeppink")
plot_start_param

Warning message:
"package 'ggplot2' was built under R version 4.0.5"
Linkoping Temperature μ_0, Ω_0, ν_0



In [22]: lmTemp\$coefficients

(Intercept): -11.9556531804222 time: 103.584049398349 I(time^2): -95.418518926651

adjusting $\mu_0, \Omega_0, \nu_0, \sigma_0^2$

The starting value results to very high value for temprature (i.e. 150°C). This is unreasonable for swedish weather. To achieve better prior we adjust the model parameter as following:

we decided to set the initial value of μ_0 to the `lmTemp$coefficients` we calculated earlier.

$\mu_0 = (-12, 103, -95)$

From the above plot we see lots of variation in the models so we decided to reduce the value of σ_0 to 0.03. This decision was made by trial and error.

we also increased the value of ν_0 to 10.

In [23]: mu.0 = c(-11, 103, -95)
omega.0 = 0.01*diag(3)
nu.0 = 10
sigma2.0 = 0.03
Ndraws = 100
Bayes.Regressor2 = matrix(, length(time), Ndraws)
ErrorTerm2 = numeric(Ndraws)
sigma2 = numeric(Ndraws)
BetaList2 = matrix(, Ndraws, 3)
colnames(BetaList2) = c("B0", "B1", "B2")

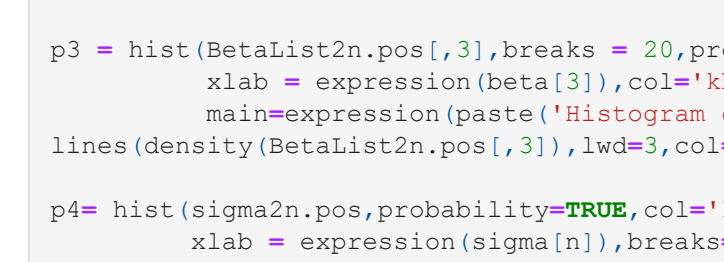
for(i in 1:Ndraws){
 sigma2[i] = sigma2.prior()
 BetaList2[i,1] = Beta.prior(sigma2[i])[1]
 BetaList2[i,2] = Beta.prior(sigma2[i])[2]
 BetaList2[i,3] = Beta.prior(sigma2[i])[3]
 ErrorTerm[i] = rnorm(1, mean = 0, sd = sqrt(sigma2))
}

for(i in 1:Ndraws){
 Bayes.Regressor2[i,1] = BetaList2[i,1] + BetaList2[i,2]*time +
 BetaList2[i,3]*(time^2) + ErrorTerm[i]
}

TempLink2 = data.frame(TempLink)
models2 = data.frame(TempLink2, Bayes.Regressor2)
plot_new_param = ggplot(models2, aes(y=temp, x = time)) +
 labs(title = expression(paste("Linkoping Temperature revised value for ",
 mu[0], " ", sigma[0], " ", Omega[0], " ", nu[0])), x = "Time", y = "Temperature") + theme_minimal()

for(i in names(models2)[-c(1,2)]){
 plot_new_param = plot_new_param +
 geom_line(aes_string(y = i), color = "blue", alpha = 0.2)
}

plot_new_param = plot_new_param +
 geom_point(aes(y = temp), alpha = 0.5, color = "deeppink")
plot_new_param



In [24]:

In [25]:

In [26]:

In [27]:

In [28]:

In [29]:

In [30]:

In [31]:

In [32]:

In [33]:

In [34]:

In [35]:

In [36]:

In [37]:

In [38]:

In [39]:

In [40]:

In [41]:

In [42]:

In [43]:

In [44]:

In [45]:

In [46]:

In [47]:

In [48]:

In [49]:

In [50]:

In [51]:

In [52]:

In [53]:

In [54]:

In [55]:

In [56]:

In [57]:

In [58]:

In [59]:

In [60]:

In [61]:

In [62]:

In [63]:

In [64]:

In [65]:

In [66]:

In [67]:

In [68]:

In [69]:

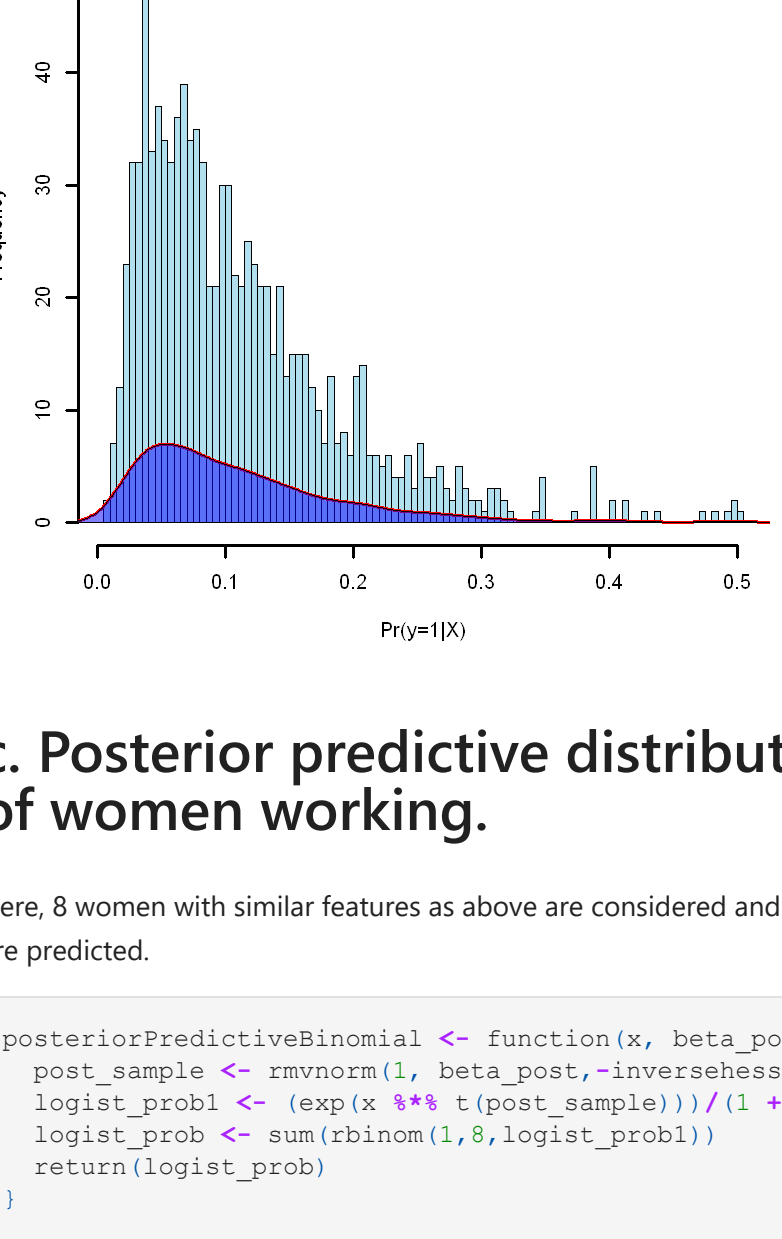
In [70]:

In [71]:

In [72]:

In [73]:

Posterior Predictive Plot



c. Posterior predictive distribution for the number of women working.

Here, 8 women with similar features as above are considered and the probabilities of how many them work are predicted.

```
In [56]: posteriorPredictiveBinomial <- function(x, beta_post, inversehessian){
  post_sample <- rmvnorm(1, beta_post, inversehessian)
  logist_prob1 <- (exp(x %*% t(post_sample)))/(1 + exp(x %*% t(post_sample)))
  logist_prob <- sum(rbinom(1,8,logist_prob1))
  return(logist_prob)
}

In [57]: test_data <- matrix(x,nrow=8, ncol=8,byrow = TRUE)
nsamples = 1000
post_predict <- c(rep(0,nsamples))
for(i in 1:nsamples){
  post_predict[i] <- posteriorPredictiveBinomial(test_data, beta_post,
  inversehessian)
}
hist(post_predict,breaks = 100,col = 'lightblue2',
  xlab = 'Number of women who works',
  main = 'Histogram of number of women who works ')
```

