

Machine Learning II

Week 1 Lecture – 9th October 2019

Imputing Missing Values, Gibbs Sampling

Topics in Machine Learning 1

Unsupervised Learning

- ▶ K-Means, K-Medoids
- ▶ Hierarchical Clustering
- ▶ Soft Clustering: EM Algorithm

Supervised Learning

- ▶ Regression Problems
 - Linear and multiple regression, ridge regression, the lasso.
- ▶ Classification Problems
 - Bayes classifier (theoretical), logistic regression, linear and quadratic discriminant analysis
- ▶ Tree Methods Regression & Classification
- ▶ Ensemble Methods (applied to tree models)

Machine Learning II

This course builds on the ideas learnt in the supervised learning part of ML1.

The main topics are

- ▶ Coping with missing values,
- ▶ Nonlinear regression methods incl. Spline Smoothing and GAMs,
- ▶ Naive Bayes Classifier,
- ▶ Support Vector Machines

The last five weeks are focussed on **Artificial Neural Networks**.

A provisional week by week semester plan is given in Moodle.

As in ML1, the methods will be applied using the software R.

Course Introduction

To pass the course you need to gain at least 45 out of 100 course marks.

The course is examined 70% by a 90 minute written exam. Everything covered in the course including computational aspects are examinable unless explicitly stated by the lecturer. Aspects covered in Machine Learning 1, that are directly relevant to Machine Learning 2, such as MSE, are also examinable.

In addition there will be a major project worth 30 Marks. The project is not compulsory, but in practise passing the the course will be difficult without handing one in.

- ▶ In the 1st exam period: 29th January 2020 at 14:15.
- ▶ In the 2nd exam period: 29th March 2019 at 10:00 (provisional).

The exam locations are still to be organised.

Information on all exam timetabling can be found at:

<https://pruefungen.beuth-hochschule.de/M-DS>

Project

- ▶ You will work in small groups analysing a medium sized data set using supervised learning.
- ▶ Methods learnt in (the first part) of *Machine Learning 2* are relevant.
- ▶ You can also use methods learnt in *Machine Learning 1* but are but not compulsory.
- ▶ Finding a suitable data set is a part of the project. You should check with the lecturer that the data are appropriate before analysing the data.
- ▶ The provisional time plan for the project is:
 - 27th November – Project details given out.
 - 11th December or before – details of your proposed data set and the students in your group to be submitted to the lecturer.
 - 10th January 2019 – Deadline for submitting your work.

Coping With Missing Data

- ▶ Introduction
- ▶ Types of missing data
- ▶ Univariate imputation methods: Fixed methods
- ▶ Univariate imputation methods: Random methods
- ▶ Multivariate imputation
- ▶ MCMC and Gibbs Sampling
- ▶ Imputation using Gibbs sampling and quick example

Introduction

There is a `Diabetes2` data set you will be using later in the course.

These data originated from the National Health and Nutrition Examination surveys (NHANES) in the USA. The full data contains 10 000 elements but there are a lot of missing values.

When selecting 13 variables, only three variables contain complete data. One of these 13 variables has 3327 missing values.

```
> Diabetes2[100:102,]
      YN Gender  Race1  BMI Age Pulse BPSysAve BPDiaAve HealthGen
100  No female Mexican 18.62 12   74      120      55      Good
101  Yes  male   White 34.04 80   60      125      56      Good
102  Yes female   Black 28.60 66  110      136      70      <NA>

      DaysPhysHlthBad DaysMentHlthBad LittleInterest Depressed
100                5                0             <NA>      <NA>
101                0                0             None      None
102               NA               NA             <NA>      <NA>
```

There is an R-Function called `na.omit()`, which returns the input vector/data frame excluding any element containing a missing value. Many data analysis programs do this internally to strip any missing values.

This is the quickest and easiest way to deal with missing values and what we have implicitly used until now.

`na.omit()` on the `Diabetes2` data results in a data set with 6492 elements, we are throwing away a third of our data.

If we were to choose to use more than 13 variables, then we would have to throw out even more elements.

A quick and easy method is usually a poor method.

Aim of imputing missing values

Imputation is the name for any method which estimating data using the data itself and is usually applied to filling-in the missing values in a data set.

The aim of imputing missing values is **not** to exactly predict the correct value.

The aim is to complete the data so that the entire data can be used, therefore increasing the model accuracy by reducing the variance.

The imputation process should not bias estimates and predictions of a machine learning method.

Types of missing data

Missing completely at random (MCAR) The probability that a missing value occurs¹ is constant. The occurrence is independent from any variable in the data

Missing at random (MAR) The probability is not constant, but can be fully explained using variables in the data set.

E.g. Women are (supposedly) less likely to provide their age than men. As long as missing values in an age variable depend only on gender, which is a variable in the data then the data are MAR.

Missing not at random (MNAR) The probability that a missing value occurs depends on non observed variables.

¹“missing value occurs” is often called “missingness”

Example of MNAR data: In a particular clinical study, patients who were given a new treatment were more likely to get headaches and drop out of the study. If the side effect information is not recorded, then the data collected at the end of the study is missing not at random (MNAR).

There is a special case of missing not at random which can be a particular problem. When the missing status of a variable depends on the value of that variable itself. E.g. Obese people are less likely to report their weight.

For MCAR or MAR data then omitting the data will not bias the results, but we will lose prediction power. The variance of our estimates will increase, and so will the mean squared error MSE.

It is usually difficult to assess whether the missing data is MCAR, MCAR or MNAR.

Univariate imputation methods: Fixed methods

Mean replacement For each variable independently, calculate the mean of the non-missing values and set the missing values equal to that mean.

E.g. `pulse` rate has 1437 missing values, the mean pulse rate of the non missing values is 73.56. Mean replacement will replace the missing values with 73.56.

This method is not good because the standard deviation of the imputed variable will decrease, as will correlations.

E.g. The standard deviation of the non missing values for `pulse` is 12.15 but after mean replacement it drops to 11.25

Missing as a category

For a **categorical variable** add a new level called missing. E.g. for `gender` define three levels: female, male and missing. This is a reasonable a method.

For a **numerical variable** replace the missing values with the mean or median value, and create a new variable called missing. This is also not a good method, as this leads to biased coefficients for the other variables.

Using logic rules

Sometimes a value can be imputed exactly or with high accuracy.

E.g. if the a person's `weight` and `BMI` is known, but `height` is missing, then it can be calculated from the definition of body mass index.

E.g. If a variable is number of pregnancies and a positive number is given then a missing value for `gender` can be imputed as female.

E.g. Only 1% of breast cancer cases are in males. This means that if there is an indication of current or previous breast cancer, it is acceptable to impute that the `gender` is female.

Univariate imputation methods: Random methods

Mean/Variance Simulation

Suppose the variable x_1 has missing values.

Compute the mean \bar{x}_1 and the standard deviation s_{x_1} of the non-missing values.

Replace the missing values with simulations from a normal $N(\bar{x}_1, s_{x_1}^2)$ distribution.

This is a simple and effective method provided x_1 is roughly normal distributed. If the distribution is very non-normal, then the variable after imputation will be closer to the normal distribution.

The disadvantage with this and all univariate methods, is that it does not consider the co-dependencies other variables.

Direct Random Sampling

Sample the non missing values with replacement, to fill in the missing values.

```
> tt
[1] 3 1 0 2 0 2 NA 2 NA 1 5 NA 0 3 4 NA 0 4 4 NA
> tt[!is.na(tt)]
[1] 3 1 0 2 0 2 2 1 5 0 3 4 0 4 4
> tt[is.na(tt)]<-sample(tt[!is.na(tt)],5,replace=TRUE)
> tt
[1] 3 1 0 2 0 2 2 2 1 1 5 2 0 3 4 3 0 4 4 2
```

This approach is good when the missingness is independent of all other variables. Compared to the *mean/variance simulation* method, direct random sampling will preserve the distributional properties of the variable, and can be used for all data types.

It is often used as a first step with multivariate methods.

Multivariate imputation for one variable

With most data there is codependency (correlation) between some variables, which can and should be incorporated into our imputation method.

E.g. If we know a person's `height`, `age` and `gender`, then we can impute the `weight` much more accurately than just sampling the non-missing values of `weight`.

We will assume from here on that the data consists entirely of continuous variables. Allowing other types of variables is more complex, but follows similar methods.

Multivariate Regression

We can take the elements with complete data and fit a regression model to predict the missing values for another variable.

The type of regression can be any regression method but is often multiple linear regression.

E.g. Assume we have complete data ($n=10\,000$) for `age`, `race`, `gender` and `BMI`. The variable `Pulse` has 1437 missing values.

We can fit a linear regression model

`Pulse ~ BMI + age + race + gender`
using the 8563 elements with complete data.

The predicted missing values should then be simulated with an error term.

If \hat{y}_i is the predicted value from the regression model, then the imputed value should be $\tilde{y} = \hat{y}_i + \epsilon_i$, where ϵ_i is a random normal simulation with the model residual variance.

As with *Mean/Variance Simulation*, simulating the error term avoids a reduction in the variance of that variable. Remember, our aim is not to get the best prediction for the missing value but to preserve the overall properties of the complete data set.

The regression method is usually used with continuous variables.

With binary variables logistic regression can be used.

For a factor variable with more than two levels use a tree model or other type of simple classifier.

Missing values with more than one variable

The *multivariate regression method* easily adapts itself to imputing multiple variables, by looping over the the variables with missing values.

Suppose x_1, \dots, x_k are variables with missing values and x_{k+1}, \dots, x_p are complete.

- ▶ Fit the regression $x_1 \sim x_{k+1} + x_{k+2} + \dots + x_p$, and impute all the missing values in x_1 including simulated error term.
- ▶ Then fit the regression $x_2 \sim x_1 + x_{k+1} + x_{k+2} + \dots + x_p$, and impute all the missing values in x_2 with simulated error term.
- ▶ Continue until all variables x_1, \dots, x_k have been completed.

A problem with this method is that if x_1 is correlated with x_2 then this is not considered when imputing x_1 , but is considered when imputing x_2

So in practice the following method is used.

- ▶ First, use a simple univariate method such as *direct random sampling* to complete each variable. Keep a record of the “positions” of the original missing values.
- ▶ Fit the regression $x_1 \sim x_2 + x_3 + \dots + x_p$, using all elements for which x_1 was not imputed.
- ▶ Re-impute all the originally missing values in x_1 including simulated error term.
- ▶ Then fit the regression $x_2 \sim x_1 + x_3 + x_4 + \dots + x_p$, etc.
- ▶ Continue until x_1, \dots, x_k have been completed.
- ▶ Because the x_1 imputations were based on some *direct random sampled* values of $x_2 \dots, x_k$, it is best to repeat this loop a couple more times.

This method is a variant of the EM-Algorithm (ML I Lecture 3: soft clustering)

Imputing missing values using Gibbs sampling

Paraphrasing Mean/Variance Simulation (Slide 14)

Replace the missing values of x_1 with simulations from a normal $N(\hat{\mu}_1, \hat{\sigma}_1^2)$ distribution.

...

The disadvantage with this and all univariate methods, is that it does not consider the co-dependencies other variables.

Gibbs sampling (GS) is a method which simulates the missing data by estimating $\hat{\mu}_1, \hat{\sigma}_1, \hat{\mu}_2, \hat{\sigma}_2, \dots, \hat{\mu}_p, \hat{\sigma}_p$ simultaneously.

In fact, in each iteration you get a slightly different value for each of these parameter estimates which reflects that the estimates are not known exactly.

What is Gibbs Sampling?

In week 4 you will learn about **Bayesian estimation** and conditional independence.

There is a wide class of simulation algorithms called **Markov Chain Monte Carlo (MCMC)** simulation which uses Bayesian estimation. These methods allow us to simulate parameter values that come from the so called **posterior distribution** of these parameters.

Simulating from the posterior distribution using MCMC is a powerful tool in Bayesian estimation.

Gibbs sampling is one of the most common MCMC algorithms. Each time a new parameter value is simulated it depends on the values of the known data, the values of the other parameters, and the other imputed values.

More details are in the appendix at the end of these notes. We will only consider Gibbs sampling as a method of imputing missing values via the R-package `MICE`.

Imputation using Gibbs Sampling in R: Quick example

R-package `mice`: Multivariate Imputation by Chained Equations.
Uses Gibbs sampling as described above.

This package easily allows for multiple realisations of the complete data set.
The data analysis (if using R) can then be repeated on each of the realisations (using one command) and an aggregate model obtained.

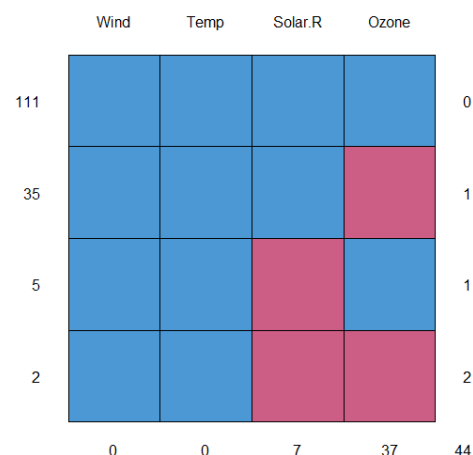
Data: `airquality` data set in R.

“Daily air quality measurements in New York, May to September 1973.”

Variables: `Ozone`, `Solar.R` (Solar radiation in Langley), `Wind` and `Temp`.
The data frame `airquality` also contains date data which we will ignore.

Inspect the pattern of missing values

```
> library(mice)
> md.pattern(data)
      Wind Temp Solar.R Ozone
111    1    1      1     1
35     1    1      1     0
5      1    1      0     1
2      1    1      0     0
      0    0      7    37 4
```



`Wind` and `Temp` are complete.

`Solar.R` has 7 missing values.

`Ozone` has 37 missing values.

Two rows have both `Solar.R` and `Ozone` missing

The function `mice` runs the Gibbs sampler.

`m=5`: obtain 5 different imputations (realisations)

`meth=norm`: the method used to obtain the conditional distributions and updates, `norm` means “Bayesian linear regression”

```
> tempData <- mice(data,m=5,maxit=50,meth='norm',seed=500)
  iter imp variable
    1   1 Ozone   Solar.R
... snip ...
```

These are the imputed values for `Ozone`. The 5 columns are the replications, and the row numbers are given.

```
> tempData$imp$Ozone
      1    2    3    4    5
5      14   18   19   19  18
10     41   45   44   32  22
25      6   19   18   19  18
... snip ...
115    21   22   32   14  13
119    77   76   50   85  66
150     9   24   23   41  13
```

To obtain a full data set with the 1st imputation.

Check using `md.pattern()`

```
> completedData <- complete(tempData,1)
> md.pattern(completedData)
/\      /\
{  '----'  }
{  O    O  }
==>  V <== No need for mice. This data set is completely observed.
\  \|\ /  /
  '-----'
```

```
      Ozone Solar.R Wind Temp
153      1       1    1    1  0
      0       0    0    0  0
```

Fit a linear regression model **to each of the realisations** and inspect the output.

```
> modelFit1 <- with(tempData, lm(Temp~ Ozone+Solar.R+Wind))
> summary(pool(modelFit1))
```

	estimate	std.error	statistic	df	p.value
(Intercept)	71.927280378	2.782518407	25.849705	107.83634	0.000000e+00
Ozone	0.173190306	0.024925656	6.948275	56.09270	2.933880e-10
Solar.R	0.009809583	0.007712214	1.271954	35.91124	2.061256e-01
Wind	-0.305981221	0.206219420	-1.483765	97.51888	1.407882e-01

Further Reading

Ashcroft, Module 6 Missing Data.

Essentials of Data Analytics and Machine Learning [Link](#)

Gelman and Hill, Chapter 25 Missing-data imputation.

Data Analysis Using Regression and Multilevel/Hierarchical Models. [Link](#)

Appendix: Background theory

The details here are not examinable, unless they come up later in another lecture.

The Bayesian approach In statistics there are two main approaches to estimating a parameter, frequentist and Bayesian.

In a very simple case we have a sample x_1, \dots, x_n and we estimate the mean of the underlying population (μ).

The frequentist approach says: we estimate μ using the sample mean, so

$$\hat{\mu} = \bar{x} = \frac{\sum x_i}{n}.$$

The Bayesian approach says: we have a very rough idea about μ , this is called our *prior knowledge*.

We assign a distribution to this prior knowledge.

Our *prior parameter* μ_0 is assigned a *prior distribution* $\mu_0 \sim \pi_0$.

Once we have obtained the data x_1, \dots, x_n we can update our information about μ to give us a **posterior distribution** $\mu_1 \sim \pi_1$ using Bayes theorem.

This posterior distribution is more accurate than the prior because it makes use of the data.

Often μ_1 is written $\mu|X$ “mu given the data”.

Bayes Theorem: from ML1 Lecture 7 (24th May)

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Let A be the parameter μ and B be the data x_1, \dots, x_n

Very roughly $P(A)$ is the prior distribution π_0 and $P(A|B)$ is the posterior distribution π_1 given the data.

MCMC

An explicit formula for the posterior is only realistic under very specific circumstances.

Very often an algorithm is used to simulate parameter values coming from the posterior distribution.

Suppose we now have many parameters expressed as a vector θ . The posterior distribution is $\pi_1(\theta|x_1, \dots, x_p)$.

If we can simulate a vector θ which comes from the posterior distribution, then we can use many simulated values to make inferences.

E.g. we simulate 1000 values of θ which includes 1000 simulations of $\theta_1 = \mu_1$, we choose the mean of these 1000 values to be our estimate $\hat{\mu}$.

There is a wide class of such simulation algorithms called **Markov Chain Monte Carlo (MCMC)** simulation. Simulating from the posterior distribution using MCMC is a powerful tool in Bayesian statistics.

The two most common MCMC algorithms are **Gibbs sampling** and the **Metropolis-Hastings Algorithm**.

The Metropolis-Hastings (MH) Algorithm is very adaptable and easy to code. The disadvantage is that it uses a rejection algorithm, and so typically twice as many iterations are needed as in Gibbs sampling.

Gibbs sampling is easier to understand at the theoretical level and is quicker than MH, but more work is needed to define the Gibbs sampler probabilities and can be more sensitive to model assumptions.

All MCMC methods apply an iterative procedure. Arbitrary starting values for θ are assigned, and new values of θ are iteratively simulated using the “Markov chain”.

After a **burn-in** period the Markov chain converges and the values of θ after each full iteration are simulations from the posterior distribution

$$\pi(\theta | x_1, \dots, x_p).$$

Gibbs sampling

If we know the conditional distribution of θ_1 given everything else, e.g. $\theta_1 | \theta_2, \dots, \theta_L, x_1, \dots, x_p$, then we can simulate a value for θ_1 assuming the other values are known.

We update θ_1 with this simulated value using the *current* values of $\theta_2, \dots, \theta_L$.

Then move on to θ_2 :

If we know the conditional distribution of θ_2 given everything else, e.g.

$\theta_2|\theta_1, \theta_3, \dots, \theta_L, x_1, \dots, x_p$, then we can simulate a value for θ_2 assuming the other values are known.

We update θ_2 with this simulated value using the *current* values of $\theta_1, \theta_3, \dots, \theta_L$.

etc.

Once we have completed updated the value of θ_L , we have finished one iteration, and have a new value for θ .

Statistical theory shows that, after the burn in period, each updated θ is a simulation from the posterior distribution $\pi(\theta|x_1, \dots, x_p)$.

Gibbs sampling to impute missing values

In the above algorithm the data x_1, \dots, x_p are known and fixed.

If x_1 contains missing values then we can also update x_1 given everything else.

If x_1 has, say, 10 values to be imputed, then we simulate 10 values from the conditional distribution $x_1|\theta_1, \dots, \theta_L, x_2, \dots, x_p$

Continue updating all variables with missing data.

$x_2|\theta_1, \dots, \theta_L, x_1, x_3, \dots, x_p$ etc.

A full iteration now consists of updating all the θ s and all the missing values.

Reminder: the aim is not exact predictions of the missing, and so we only require one complete simulation once the burn in period is complete. But by running the algorithm a few more iterations, we can obtain several different fully imputed data sets.