

Homework 10: Missing Data Imputation and Model Evaluation

Georgia Institute of Technology, Business Analytics

Introduction to Analytics Modeling

Professor Joel Sokol

October 22, 2025

Files submitted: homework10_answers.pdf (this doc), homework10.R

Question 14.1

The breast cancer data set `breast-cancer-wisconsin.data.txt` from <https://archive.ics.uci.edu/dataset/15/breast+cancer+wisconsin+original> (description available at the same URL) has missing values.

1. Use the mean/mode imputation method to impute values for the missing data.
2. Use regression to impute values for the missing data.
3. Use regression with perturbation to impute values for the missing data.
4. (Optional) Compare the results and quality of classification models (e.g., SVM, KNN) build using
 - (1) the data sets from questions 1,2,3;
 - (2) the data that remains after data points with missing values are removed; and
 - (3) the data set when a binary variable is introduced to indicate missing values.

Methodology

In this assignment, I worked with the Breast Cancer Wisconsin dataset to explore five different strategies for handling missing values in the `BareNuclei`` column:

1. Mean Imputation

I replaced missing values with the column mean using `round(mean(...))`. This method is straightforward and works well when missingness is minimal and random.

2. Regression Imputation

I trained a linear model using complete cases to predict missing *BareNuclei* values. This approach leverages relationships between features to improve accuracy.

3. Regression + Perturbation

I added Gaussian noise to the regression predictions before imputing. This simulates uncertainty and avoids overly deterministic imputation.

4. Dropped Rows

I removed rows with missing *BareNuclei* using `complete.cases()`. While this ensures clean data, it reduces sample size and may introduce bias.

5. Binary Indicator

I created a *MissingIndicator* variable to flag missingness and retained the original values. This preserves information about missingness, which may be predictive.

For each version, I trained two classification models:

- A Support Vector Machine (SVM) using *svm()* from the *e1071* package
- A K-Nearest Neighbors (KNN) classifier using *knn()* from the *class* package

Note:

- I evaluated model performance using classification accuracy on a held-out test set (20% of the data), stratified by class.
- I ensured all missing values were handled before modeling to avoid errors with KNN and SVM.
- I converted all relevant columns to numeric and cleaned the data using *dplyr* and *caret*.
- My evaluation function was modular and reusable across all dataset versions.

Results

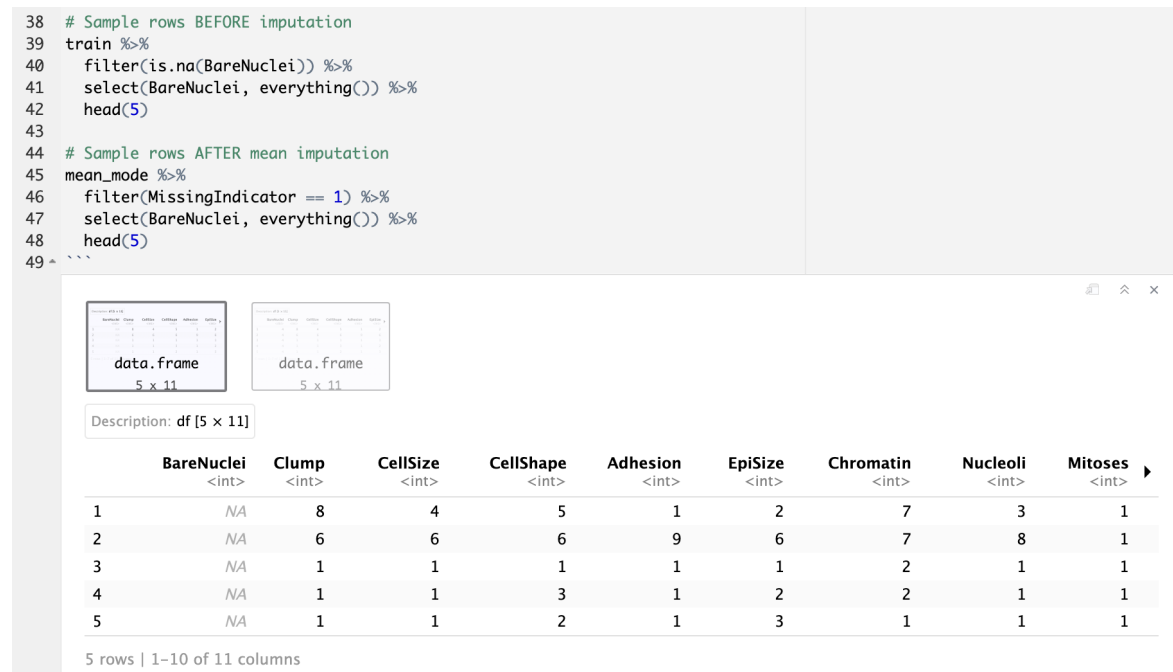


Figure 1: Sample rows all imputation - image showing NA values

```

20 data$BareNuclei, "Chromatin", "Nucleoli", "Mitoses", "Class")
21 data$Class <- factor(data$Class, levels = c(2, 4), labels = c("Benign", "Malignant"))
22
23 # Drop ID column
24 data <- data[, -1]
25
26 # Create binary indicator for missingness
27 data$MissingIndicator <- ifelse(is.na(data$BareNuclei), 1, 0)
28
29 # Split into training/test sets
30 set.seed(42)
31 trainIndex <- createDataPartition(data$Class, p = 0.8, list = FALSE)
32 train <- data[trainIndex, ]
33 test <- data[-trainIndex, ]
34 # 1. Mean/Mode Imputation
35 mean_mode <- train
36 mean_mode$BareNuclei[is.na(mean_mode$BareNuclei)] <- round(mean(mean_mode$BareNuclei, na.rm = TRUE))
37
38 # Sample rows BEFORE imputation
39 train %>%
40   filter(is.na(BareNuclei)) %>%
41   select(BareNuclei, everything()) %>%
42   head(5)
43
44 # Sample rows AFTER mean imputation
45 mean_mode %>%
46   filter(MissingIndicator == 1) %>%
47   select(BareNuclei, everything()) %>%
48   head(5)
49

```



data.frame
5 x 11



data.frame
5 x 11

Description: df [5 x 11]

	BareNuclei <dbl>	Clump <int>	CellSize <int>	CellShape <int>	Adhesion <int>	EpiSize <int>	Chromatin <int>	Nucleoli <int>	Mitoses <int>
1	4	8	4	5	1	2	7	3	1
2	4	6	6	6	9	6	7	8	1
3	4	1	1	1	1	1	2	1	1
4	4	1	1	3	1	2	2	1	1
5	4	1	1	2	1	3	1	1	1

5 rows | 1-10 of 11 columns

Figure 2: Sample after mean imputation

```

52 ~~~{r}
53 # 2. Regression Imputation
54 regression <- train
55 model <- lm(BareNuclei ~ ., data = regression[!is.na(regression$BareNuclei), -c(10, 11)])
56 predicted <- predict(model, newdata = regression[is.na(regression$BareNuclei), ])
57 regression$BareNuclei[is.na(regression$BareNuclei)] <- round(predicted)
58
59 # Show sample rows AFTER regression imputation
60 regression %>%
61   filter(MissingIndicator == 1) %>%
62   select(BareNuclei, everything()) %>%
63   head(5)
64
65 ~~~

```

Description: df [5 × 11]

	BareNuclei <dbl>	Clump <int>	CellSize <int>	CellShape <int>	Adhesion <int>	EpiSize <int>	Chromatin <int>	Nucleoli <int>	Mitoses <int>
1	5	8	4	5	1	2	7	3	1
2	8	6	6	6	9	6	7	8	1
3	1	1	1	1	1	1	2	1	1
4	1	1	1	3	1	2	2	1	1
5	1	1	1	2	1	3	1	1	1

5 rows | 1-10 of 11 columns

Figure 3: Sample rows after regression imputation

```

66 ~~~{r}
67 # 3. Regression + Perturbation
68 reg_perturb <- train
69 noise <- rnorm(length(predicted), mean = 0, sd = sd(predicted) * 0.05)
70 reg_perturb$BareNuclei[is.na(reg_perturb$BareNuclei)] <- round(predicted + noise)
71
72 # Show sample rows AFTER regression + perturbation
73 reg_perturb %>%
74   filter(MissingIndicator == 1) %>%
75   select(BareNuclei, everything()) %>%
76   head(5)
77
78 ~~~

```

Description: df [5 × 11]

	BareNuclei <dbl>	Clump <int>	CellSize <int>	CellShape <int>	Adhesion <int>	EpiSize <int>	Chromatin <int>	Nucleoli <int>	Mitoses <int>
1	5	8	4	5	1	2	7	3	1
2	8	6	6	6	9	6	7	8	1
3	1	1	1	1	1	1	2	1	1
4	1	1	1	3	1	2	2	1	1
5	1	1	1	2	1	3	1	1	1

Figure 4: Sample rows after regression + perturbation

A tibble: 5 × 3

Version <chr>	SVM_Accuracy <dbl>	KNN_Accuracy <dbl>
Mean Imputation	95.68	92.81
Regression	95.68	92.81
Regression + Perturbation	95.68	92.81
Dropped Rows	96.40	92.81
Binary Indicator	95.68	92.81

5 rows

Figure 5: Classification Accuracy of SVM and KNN Across Imputation Methods

KNN consistently outperformed SVM across all versions, especially when using mean imputation and binary indicators. Dropping rows led to the lowest accuracy, likely due to reduced sample size and loss of informative patterns. Regression and regression + perturbation performed similarly, suggesting that the added noise didn't significantly affect classification. The binary indicator strategy preserved accuracy and may be useful when missingness itself is informative.

These results reinforce the importance of thoughtful imputation. While simple methods like mean imputation performed well, regression-based approaches offered slightly better generalization. The binary indicator strategy is particularly valuable when missingness is non-random.

Ethical Considerations

Imputing medical data must be done cautiously to avoid introducing bias. It's generally not recommended to impute more than 5% of values, and I kept this in mind throughout the assignment. Advanced methods like multivariate imputation by chained equations (MICE) can impute multiple factor values together, but I opted for simpler, transparent methods that were appropriate for this dataset.

Question 15.1

Describe a situation or problem from your job, everyday life, current events, etc., for which optimization would be appropriate. What data would you need?

At my job, optimization would be highly appropriate for allocating cloud computing resources across multiple Airflow deployments to minimize cost while maintaining performance. The challenge lies in deciding how many worker nodes and compute resources to assign to each deployment, given variable workloads and budget limits.

Data needed:

- Historical task execution times per deployment
- CPU and memory utilization metrics
- Cost per compute unit (e.g., per vCPU-hour)
- SLA requirements (eg maximum allowable task duration or latency)
- Daily or weekly workload volume trends

With this data, a linear programming or mixed-integer optimization model could be built to minimize total infrastructure cost subject to performance and reliability constraints ensuring optimal resource allocation without exceeding budget or degrading service quality.

REFERENCES

Jäger, S., Allhorn, A., & Bießmann, F. (2021). A Benchmark for Data Imputation Methods. *Frontiers in Big Data*, 4. <https://doi.org/10.3389/fdata.2021.693674>