Applied ML PreProcessing of Data







1. S. Mar. 2. 17 - 183





Troels C. Petersen (NBI)



"Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!"

The problem of real data!

In an ideal world, we would:

Copy data to your area and start training on it!

The problem of real data!

In an ideal world, we would:

Copy data to your area and start training on it!

In the real world, we need to:

- Collecting data from detectors, databases, spreadsheets, APIs, etc. ensuring that it contains the relevant parts in an understood way.
- Clean the data by removing or correcting missing values, outliers, or inconsistencies.
- Transform the data through normalisation and encoding to make it compatible with machine learning algorithms.
- Potentially, reduce the data's complexity without losing the information through dimensionality reduction.

The following are about these steps.

So far, we have looked at "perfect" data, i.e. data without any flaws in it. However, real world datasets are hardly ever "perfect", but contains flaws that makes preprocessing imperative.

Effects may be (non-exhaustive list):

- NaN-values and "Non-values" (i.e. -9999)
- Wild outliers (i.e. values far outside the typical range)
- Shifts in distributions (i.e. part of data having a different mean/width/etc.)
- Mixture of types (i.e. numerical and text, from something humans filled out)

It is also important to consider, if entries are missing...

- 1. Randomly (in which case there should be no bias from omitting) or
- 2. Following some pattern (in which case there could be problems!).

In case of NaN values, we might simply decide to drop the variable column or entry row, requiring that all variables/entries have reasonable values.

Alternatively, we might insert "**imputed**" values instead, saving statistics.

Imputations (dealing with missing values)

Imagine that you get the following data, which contains (many) NaN values. How would you want to treat a dataset like this? Discuss locally and classwide.

ID	City	Degree	Age	Salary	Married ?
1	Lisbon	NaN	25	45,000	0
2	Berlin	Bachelor	25	NaN	1
3	Lisbon	NaN	30	NaN	1
4	Lisbon	Bachelor	30	NaN	1
5	Berlin	Bachelor	18	NaN	0
6	Lisbon	Bachelor	NaN	NaN	0
7	Berlin	Masters	30	NaN	1
8	Berlin	No Degree	NaN	NaN	0
9	Berlin	Masters	25	NaN	1
10	Madrid	Masters	25	NaN	1

Medium Article, "Working with Missing Values in a Dataset" by Okoh Anita 6

Removing Variables/Columns

A "simple" way to get rid of NaN-values is to drop the column(s) containing NaN values. This works well, when these are few and with a high NaN-fraction.

ID	City	Salary	Married ?		ID	City	Married ?
1	Lisbon	45,000	0		1	Lisbon	0
2	Berlin	NaN	1	~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~	2	Berlin	1
3	Lisbon	NaN	1	-	3	Lisbon	1
4	Lisbon	NaN	1	·	4	Lisbon	1
5	Berlin	NaN	0		5	Berlin	0
6	Lisbon	NaN	0		6	Lisbon	0
7	Berlin	NaN	1		7	Berlin	1
8	Berlin	NaN	0		8	Berlin	0
9	Berlin	NaN	1		9	Berlin	1
10	Madrid	NaN	1		10	Madrid	1

Removing Entries/Rows

Alternatively, one might instead decide to remove rows with many NaNs in them, arguing that these are incomplete entries. This works well, when these seem to be Random NaNs, and when the entries with a high NaN-fraction are few.

_					
10	D	City	Degree	Age	Married ?
1	I	Lisbon	NaN	25	0
2	2	Berlin	Bachelor	25	1
3	3	Lisbon	NaN	30	1
4	1	Lisbon	Bachelor	30	1
5	5	Berlin	Bachelor	18	0
6	6	Lisbon	Bachelor	NaN	0
7	7	Berlin	Masters	30	1
8	3	Berlin	No Degree	NaN	0
9)	Berlin	Masters	25	1
1	10	Madrid	Masters	25	1
· · · · · · · · · · · · · · · · · · ·					

So, in case of NaN values, we might simply decide to drop the variable column or entry row, requiring that all variables/entries have reasonable values.



Alternatively, we might insert "imputed" values instead, saving statistics:

Oriç	ginal		data.f	illna(0)	data.fillna(r	nethod = 'bfill')) da	ata.fillna(m	ethod
One	Two		One	Two	One	Тwo		One	Tw
0	2		0	2	0	2		0	2
1	3		1	3	1	3		1	3
NaN	0		0	0	2	0		1 🕈	0
2	1	1	2	1	2	1		2	1

Imputing with the mode

Instead of removing rows with NaN-values in them, one can decide to rather "impute" values into these, e.g. using the mode (most frequently occurring).

ID	City	Degree	Married ?	
1	Lisbon	NaN	0	
2	Berlin	Bachelor	1	
3	Lisbon	NaN	1	
4	Lisbon	Bachelor	1	
5	Berlin	Bachelor	0	
6	Lisbon	Bachelor	0	
7	Berlin	Masters	1	
8	Berlin	No Degree	0	
9	Berlin	Masters	1	
10	Madrid	Masters	1	

Most frequent Degree = 'Bachelor'

ID	City	Degree	Married ?
1	Lisbon	Bachelor	0
2	Berlin	Bachelor	1
3	Lisbon	Bachelor	1
4	Lisbon	Bachelor	1
5	Berlin	Bachelor	0
6	Lisbon	Bachelor	0
7	Berlin	Masters	1
8	Berlin	No Degree	0
9	Berlin	Masters	1
10	Madrid	Masters	1

Imputing with the mean

Alternatively, one may impute values using the mean (or mode) of that column. This has the advantage, that variables/events are not thrown away.

ton 1				Avera	ge_Age = 26.0	0		
	ID	City	Age	Married ?	ID	City	Age	Married ?
-	1	Lisbon	25	0	1	Lisbon	25	0
_	2	Berlin	25	1	2	Berlin	25	1
	3	Lisbon	30	1	3	Lisbon	30	1
1	4	Lisbon	30	1	4	Lisbon	30	1
0	5	Berlin	18	0	5	Berlin	18	0
	6	Lisbon	NaN	0	6	Lisbon	26	0
1	7	Berlin	30	1	7	Berlin	30	1
	8	Berlin	NaN	0	8	Berlin	26	0
	9	Berlin	25	1	9	Berlin	25	1
	10	Madrid	25	1	10	Madrid	25	1

Alternatives when imputing

Instead of using a fixed number for inputing, one might pick an age randomly, but at the same time "report" (in a new column) that it is an imputed value.

ID	City	Age	Married ?
1	Lisbon	25	0
2	Berlin	25	1
3	Lisbon	30	1
4	Lisbon	30	1
5	Berlin	18	0
6	Lisbon	NaN	0
7	Berlin	30	1
8	Berlin	NaN	0
9	Berlin	25	1
10	Madrid	25	1

The two randomly picked Ages are 25 & 30

ID	City	Age	Age-Imputed	Married ?
1	Lisbon	25	25	0
2	Berlin	25	25	1
3	Lisbon	30	30	1
4	Lisbon	30	30	1
5	Berlin	18	18	0
6	Lisbon	NaN	25	0
7	Berlin	30	30	1
8	Berlin	NaN	30	0
9	Berlin	25	25	1
10	Madrid	25	25	1

Alternatively, one might decide to give the NaN-values a specific value (here -1), which is different from all other values.

			Ar	rary Value	e = -1		
ID	City	Age	Married ?	ID	City	Age	Married '
1	Lisbon	25	0	1	Lisbon	25	0
2	Berlin	25	1	2	Berlin	25	1
3	Lisbon	30	1	3	Lisbon	30	1
4	Lisbon	30	1	4	Lisbon	30	1
5	Berlin	18	0	5	Berlin	18	0
6	Lisbon	NaN	0	6	Lisbon	-1	0
7	Berlin	30	1	7	Berlin	30	1
8	Berlin	NaN	0	8	Berlin	-1	0
9	Berlin	25	1	9	Berlin	25	1
10	Madrid	25	1	10	Madrid	25	1

For the text case, one might again simply consider making a new column, which indicates if the value given is imputed (no matter what method was used).

ID	City	Degree	Married ?
1	Lisbon	NaN	0
2	Berlin	Bachelor	1
3	Lisbon	NaN	1
4	Lisbon	Bachelor	1
5	Berlin	Bachelor	0
6	Lisbon	Bachelor	0
7	Berlin	Masters	1
8	Berlin	No Degree	0
9	Berlin	Masters	1
10	Madrid	Masters	1

Missing = 1 & Not Missing = 0

ID	City	Degree	Missing Indicator	Married ?
1	Lisbon	Bachelor	1	0
2	Berlin	Bachelor	0	1
3	Lisbon	Bachelor	1	1
4	Lisbon	Bachelor	0	1
5	Berlin	Bachelor	0	0
6	Lisbon	Bachelor	0	0
7	Berlin	Masters	0	1
8	Berlin	No Degree	0	0
9	Berlin	Masters	0	1
10	Madrid	Masters	0	1

Imputations methods

Statistical imputation:

Use regression models to predict the missing values based on other features. Use probability distributions to impute missing values.

Unsupervised imputation:

Use mean, median, or mode value of distribution. Use k-NN, i.e. average of neighbours in the k most similar rows. Use AutoEncoder to learn a compressed representation of the data.

Other imputation methods:

Arbitrary Value: Replace missing values with arbitrary value, such as -1 or 0. Multivariate: Impute missing value based on relationship with other features.

Whichever method you use, be sure to note the number and fraction of imputations and be clear about the method in your thesis/publications.

It also serves you well to try different methods to see the impact of (and thus systematic uncertainty from) your choice.

NaN-values tend to correlate

It is often seen, that several variables have the same source, and thus their NaN occurrence might be correlated with each other.

This can be tested by substituting 0's for numerical values and 1's for NaN values. By considering the correlation matrix of these substitute 0/1 values, one gets a pretty clear picture.

Typically, some entries are 100% correlated, as the source of these variables is shared.

Based on this information, one can better decide how to deal with these variables.

× × × 11 × × × 26 × 11 × × 21 × × 21 × × 21 × × 21 × × 21 × 21
x 13 1 0.27 0.27 0.4 0.1 0.12 0.13 0.06 0.15 0.16 0.09 0.47 0.07 0.033 0.51 0.48 0.38 0.13 0.09 4.01 0.09 0.22 0.39 0.49 0.49
x 16 5.27 1 5.35 5.34 525 5.11 554 5 617 5.56 5.12 5.56 5.12 5.14 5.33 5.36 5.3 5.46 515 5.11 5.11 5.1 5.21 5.15 5.85
X 4 5.27 5.35 1 5.5 5.16 4.1 5.67 5.65 5.12 5.67 5.66 5.85 4.42 4.4 4.36 4.4 5.5 4.13 4.4 5.6 5.16 5.16
x 10 0.4 0.34 0.5 1 0.19 0.13 0.12 0.1 0.11 0.04 0 0.16 0.54 0.39 0.33 0.48 0.32 0.6 0.17 0 0.66 0.17 0.0 0.1 0.16
x 8 01 000 0.16 0.11 1 0.18 0.34 0.17 0.15 0.22 0.38 0.32 0.31 0.2 100 0.1 0.5 0.19 0.33 0.22 0.32 0.19 0.3 0.1
x 15 0.12 0.11 0.1 0.1 0.1 0.10 1 0.44 0.35 0.4 0.46 0.57 0.63 0.6 0.01 0.1 0.15 0.5 0.4 0.4 0.55 0.58 0.41 0.49 0.33 0.27
x 19 0.13 046 647 012 0.34 0.44 1 0.37 0.34 0.46 0.54 0.64 0.6 406 0.09 0.09 0.07 0 0.37 0.57 0.5 0.02 0.38 0.41 0.33 0.25
x_17 0.0 4.07 6.05 4.0 1.17 0.35 0.37 1 0.39 0.48 0.57 0.59 0.59 6.6 6.6 6.1 4.4 -0.47 -0.48 -0.52 -0.57 0.24 -0.47 0.21 0.26
x 23 0.15 0.65 0.12 0.11 0.15 0.4 0.34 0.39 1 0.52 0.57 0.6 0.58 0.11 0 0.13 0.1 0.13 0.1 0.45 0.52 0.62 0.4 0.49 0.25 0.2
X_7 0.16 0.19 0.07 0.9 0.22 0.46 0.46 0.48 0.52 1 0.72 0.77 0.77 0.9 0.6 0.9 0.9 0.9 0.51 0.61 0.68 0.77 0.48 0.52 0.36 0.38
x 12 50 565 505 0 338 0.57 0.54 0.57 0.57 0.72 1 0.92 0.88 566 81 0 0.64-0.77-0.79 0.9 0.49-0.69-0.44-0.41
X 6 5.07 5.12 5.11 5.59 5.32 5.53 5.64 5.59 5.6 5.77 5.92 1 5.97 5.5 5.59 5.1 5.99 5.1 5.52 5.84 5.87 5.97 5.54 5.72 5.46 5.41
x 18 60 0 14 09 0.31 0.6 0.6 0.59 0.58 0.77 0.88 0.97 1 10 605 0.65 0.82 0.84 0.93 0.55 0.68 0.42 0.38
x 20 -0.33 -0.33 -0.42 -0.39 -0.2 -0.08 ECK 2.08 -0.11 ECK -0.09 1 0.34 0.33 0.27 ECK -0.04 -0.04 0.08 -0.11 -0.15
x 25 -0.51-0.36 -0.4 -0.33 FU -0.11 -0.01 FM -0.01 FM -0.01 FM -0.34 1 0.48 0.51 FU -0.12 FM -0.14 FF4 -0.19 -0.1
X_3 -0.48 -0.3 -0.36 -0.48 -0.1 -0.65 -0.1 -0.1 -0.13 -0.1 -0.15 -
X 9 -0.38 -0.46 -0.4 -0.32 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
X_1 0.13 000 0.0 6.08 0.19 0.4 0.37 0.47 0.31 0.51 0.64 0.62 0.63 0.0 0.1 0.57 0.59 0.61 0.35 0.42 0.23 0.18
X 5 00 0 11 0 11 0 11 0 11 0 33 0.4 0 57 0 48 0 45 0 61 0 77 0 84 0 82 0 12 0 9 00 0 57 1 0 79 0 8 0 48 0 56 0 39 0 3
x 22 36 6.11 3 328-0.55 4.5 4.5 4.52 4.52 4.68 4.79 4.87 4.84 14 146 4.07 146 4.07 14 0.59 0.79 1 0.82 0.48 0.67 0.4 0.37
x 24
x_11 -0.22 +0.21 -0.14 -0.17 -0.19 -0.41 -0.38 +0.24 -0.4 -0.48 -0.49 -0.54 -0.53 +0.49 +0.14 +0.14 +0.13 +0.15 +0.48 +0.48 +0.53 +1 +0.48 +0.37 +0.28
x 14 55 4.18 5 5 4.38 4.3 4.3 4.49 0.41 0.47 0.49 0.52 0.69 0.72 0.68 0.01 5 5 5 5 5 5 5 5 5 5 5 5 7 0.71 0.48 1 0.34 0.25
x 2 35 5 60 61 61 61 63 030 030 021 025 030 044 046 042 011 019 066 05 523 039 0.4 0.43 0.37 0.34 1 0.31
x 21 0.08 10 0.06 0.16 0.06 0.27 0.25 0.26 0.2 0.38 0.41 0.41 0.38 0.15 0.1 0.13 0.1 0.18 0.33 0.37 0.39 0.28 0.25 0.31 1
1 08 06 04 02 0 02 04 06 08

How to deal with outliers?

Sometimes, (a few?) entries take on extreme values, which ruin either the NN performance, or the transformation applied first (and then the performance). How to deal with that?

How to deal with outliers?

Sometimes, (a few?) entries take on extreme values, which ruin either the NN performance, or the transformation applied first (and then the performance). How to deal with that?



Mixture of types

At other times, types may be mixed. An example could be from a form, where people are asked "How long time did it last?". The answers:

- 2 minutes,
- 4.5h,
- about a second,
- 4:07:32,
- donno,
- all day,
- between four and 5 min.

In that case, there is nothing else to do, than to run through the bitter (large) number of unique cases, and get more than 95% working, possibly leaving the rest.

Inconsistent data

Sometimes (though more rarely), data is inconsistent. This may happen in several ways, some of which are:

- Several columns in your data, which should contain the same information, but doesn't. This often happens when you have merged (same?) data in different formats/versions.
 Example: Columns "Age" and "Days since birth" are not consistant.
- The data format may have changed during data collection, leaving a change of values at some point. Example: Column "Date" content is "2005-11-06" and later "09-08-2021".
- 3. Logical inconsistencies, where one value renders other values inconsistent. Example: Column "Car Owner" is "No", but "Own car value" is "10000\$".

There are no single recipes for dealing with these inconsistencies, except **patience**, time, dedication, and acceptance of the imperfect.

Transformations (dealing with non-unit values)

Transformations

Once we have values for all parts of data, we want to normalise it to unity, i.e. mean=0 and sigma=1.

But even perfectly good data may naturally have weird distributions that may need special attention first.

For this reason, one may want to transform the values by some power, log, etc.



Transformations

There are several ways to normalise data. In principle we want data to become Unit Gaussian, but this is not really needed, as less can do the job:

Quantile: Give values in [0,1] according to rank of values (place when sorted).

Min-Max Scaling: $z = \frac{x - min(x)}{max(x) - min(x)}$

Standardisation:

Unit Vector:

Box-Cox:

 $z = \frac{x}{||x||}$

of z most Gaussian. Implemented in SciPy.

 $z^{(\lambda)} = \begin{cases} \frac{y^{\lambda} - 1}{\lambda}, & \text{if } \lambda \neq 0\\ \ln(y), & \text{if } \lambda = 0 \end{cases}$ Input values y need to be positive. The value of λ is chosen to make the distribution

 $z = \frac{x - x}{std(x)}$

Conclusions

No matter what you plan to do with data, my first advice is always:

Print & Plot

This is your first assurance, that you even remotely know what the data contains, and your first guard against nasty surprises.

Also, working with others (from know-nothings to domain experts) you will be required to show the input, and assuring that it is valid and makes sense.

Remember also to do so in your final projects...