

# Data Pre-processing (Feature Engineering)

## Applied Machine Learning for Educational Data Science

true

08/23/2021

## Contents

<b>1. Scales of Measurement and Types of Variables</b>	<b>1</b>
<b>2. Processing Categorical Variables</b>	<b>2</b>
2.1 One-hot encoding (Dummy Variables) . . . . .	2
2.2. Label encoding . . . . .	4
2.3. Polynomial Contrasts . . . . .	5
<b>3. Processing Day and Time Data</b>	<b>5</b>
<b>4. Processing Continuous Variables</b>	<b>5</b>
4.1 Centering and Scaling . . . . .	5
4.2 Transformations . . . . .	5
4.3 Basis Expansions and Splines . . . . .	5
4.4 Linear Projections . . . . .	5
<b>Handling Missing Data</b>	<b>5</b>
<b>Data Leakage</b>	<b>5</b>
<b>Processing Text Data</b>	<b>5</b>

[Updated: Tue, Aug 24, 2021 - 18:34:38 ]

## 1. Scales of Measurement and Types of Variables

It is important to understand the nature of variables and how they were measured and represented in a dataset. In social sciences, in particular psychology, there is a methodological consensus about the framework provided by [Stevens \(1946\)](#), also see [Michell \(2002\)](#) for an in-depth discussion. According to Stevens' definition, there are four levels of measurement: nominal, ordinal, interval, and ratio. Whether a variable is considered having a nominal, ordinal, interval, or ratio scale depends on the character of the empirical operations performed while constructing the variable.

- **Nominal scale:** Variables with a nominal scale cannot be meaningfully added, subtracted, divided, or multiplied. Also, there is no hierarchical order among the assigned values. Most variables that contains labels for individual observations can be considered as nominal, e.g., hair color, city, state, ethnicity.
- **Ordinal scale:** Variables with an ordinal scale also represent labels; however, there is a meaningful hierarchy among the assigned values. For instance, if a variable is coded as Low, Medium, and High, they are simply labels but we know that High represents something more than Medium, and Medium represents something higher than Low ( $\text{High} > \text{Medium} > \text{Low}$ ). On the other side, the distance between the assigned values do not necessarily represent the same amount of difference. Other examples of variables that can be considered as ordinal are letter grades (A-F), scores from likert type items (Strongly agree, agree, disagree, strongly disagree), education status (high school, college, master's, PhD), cancer stage (stage1, stage2, stage3), order of finish in a competition (1st, 2nd, 3rd).
- **Interval scale:** Variables with an ordinal scale represents quantities with equal measurement units but they don't have an absolute zero point. For instance, a typical example of an interval scale is temperature measured on the Fahrenheit scale. The difference between 20F and 30F is the same difference as the difference between 60F and 70F. However, 0F does not indicate no heat.
- **Ratio scale:** Variables with a ratio scale represents quantities with equal measurement units and have an absolute zero. Due to the nature of the existence of absolute zero point that represent 'nothing', ratio of measurements are also meaningful. Typical examples are height, mass, distance, length.

Below table provides a summary of properties for each scale.

	Indicating Difference	Indicating Direction of Difference	Indicating Amount of Difference	Has absolute zero
Nominal	X			
Ordinal	X	X		
Interval	X	X	X	
Ratio	X	X	X	X

In this class, we classify the variables in two types: **Categorical** and **Continuous**. The variables with a nominal or ordinal scale are considered as **Categorical** and the variables with an interval or ratio scale are considered as **Continuous**.

## 2. Processing Categorical Variables

When there are categorical predictors in a dataset, it is important to translate them into numerical codes. When encoding categorical predictors, we try to preserve as much information as possible from its labels. Therefore, different strategies may be used for categorical variables with different ordinal scales.

### 2.1 One-hot encoding (Dummy Variables)

A dummy variable is a synthetic variable with two outcomes (0 and 1) to represent a group membership. When there is a nominal variable with  $N$  levels, it is typical to create  $N$  dummy variables to represent the information in the nominal variable. Each dummy variable represents a membership to one of the levels in the nominal variable. These dummy variables can be used as features in predictive models.

In its simplest case, consider variable **Race** in the Recidivism dataset with two levels: Black and White. We can create two dummy variables such that the first dummy variable represents whether or not an individual is Black and the second dummy variable represents whether or not the individual is White.

	Dummy Variable 1	Dummy Variable 2
Black	1	0
White	0	1

```

recidivism <- read.csv('https://raw.githubusercontent.com/uo-datasci-specialization/c4-ml-fall-2021/main/recidivism.csv')

table(recidivism$Race)

recidivism$black <- ifelse(recidivism$Race=='BLACK',1,0)
recidivism$white <- ifelse(recidivism$Race=='WHITE',1,0)

head(recidivism[,c('Race','black','white')])

table(recidivism$black)
table(recidivism$white)

```

Let's consider another example from the Recidivism dataset. Variable **Prison\_Offense** has five categories: Violent/Sex, Violent/Non-Sex, Property, Drug, Other. We can create five dummy variables using the following coding scheme.

	Dummy Variable 1	Dummy Variable 2	Dummy Variable 3	Dummy Variable 4	Dummy Variable 5
Violent/Sex	1	0	0	0	0
Violent/Non-Sex	0	1	0	0	0
Property	0	0	1	0	0
Drug	0	0	0	1	0
Other	0	0	0	0	1

Note that **Prison\_Offense** is missing for a number of observations. You can fill-in the missing values prior to creating dummy variables using one of the methods we will discuss later. Alternatively, we can define Missing as the sixth category to preserve that information.

	Dummy Variable 1	Dummy Variable 2	Dummy Variable 3	Dummy Variable 4	Dummy Variable 5	Dummy Variable 6
Violent/Sex	1	0	0	0	0	0
Violent/Non-Sex	0	1	0	0	0	0
Property	0	0	1	0	0	0
Drug	0	0	0	1	0	0
Other	0	0	0	0	1	0
Missing	0	0	0	0	0	1

```

table(recidivism$Prison_Offense)
names(table(recidivism$Prison_Offense))

recidivism$off_viosex <- ifelse(recidivism$Prison_Offense=='Violent/Sex',1,0)
recidivism$off_vionosex <- ifelse(recidivism$Prison_Offense=='Violent/Non-Sex',1,0)
recidivism$off_property <- ifelse(recidivism$Prison_Offense=='Property',1,0)

```

```

recidivism$off_drug      <- ifelse(recidivism$Prison_Offense=='Drug',1,0)
recidivism$off_other     <- ifelse(recidivism$Prison_Offense=='Other',1,0)
recidivism$off_missing  <- ifelse(recidivism$Prison_Offense=='',1,0)

head(recidivism[,c('Prison_Offense','off_viosex','off_vionosex','off_property','off_drug','off_other',

```

In some cases, when you have geographical location with a reasonable number of categories (e.g., counties or cities in a state, schools in a district), you can also create dummy variables to represent this information. In our case, the Recidivism dataset has a variable called **Residence\_PUMA** indicating [Public Use Microdata Area \(PUMA\)](#) for the residence address at the time individual was released. There is a total of 25 unique codes (1-25) for this variable; however, these numbers are just labels. So, one can create 25 different dummy variables to represent 25 different PUMAs.

---

## NOTE

When you fit a typical regression model without regularization using ordinary least-squares (OLS), a typical practice is to drop a dummy variable for one of the levels. So, for instance, if there are  $N$  levels for a nominal variable, you only have to create  $(N-1)$  dummy variables, as the  $N$ th one has redundant information. The information regarding to the excluded category is represented in the intercept term. It creates a problem when you put all  $N$  dummy variables into the model, because the OLS procedure tries to invert a singular matrix and you will likely get an error message.

On the other hand, this is not an issue when you fit a regularized regression model, which will be the case in this class. Therefore, you do not need to drop one of the dummy variables and can include all of them in the analysis. In fact, it may be beneficial to keep the dummy variables for all categories in the model when regularization is used in regression. Otherwise, the model may produce different predictions depending on which category is excluded.

---

## 2.2. Label encoding

When the variable of interest is ordinal and there is a hierarchy among the levels, we can still use one-hot encoding to create a set of dummy variables to represent the information in the ordinal variable. However, dummy variables will not provide any information regarding the hierarchy among categories.

For instance, consider the variable **Age\_At\_Release** in the Recidivism dataset. It is coded as 7 different age intervals in the dataset: 18-22, 23-27, 28-32, 33-37, 38-42, 43-47, 48 or older. One can create 7 dummy variables to represent each category in this variable. Alternatively, we can assign a numeric variable to each category that may represent the information in these categories. For instance, we can assign numbers from 1 to 7, respectively. Or, we can choose the midpoint of each interval to represent each category (e.g., 20,25,31,35,40,45,60).

```

require(car)
table(recidivism$Age_at_Release)

#?car::recode
recidivism$age <- recode(recidivism$Age_at_Release,
                        recodes = "'18-22' = 20;
                                '23-27' = 25;
                                '28-32' = 30;

```

```
'33-37' = 35;  
'38-42' = 40;  
'43-47' = 45;  
'48 or older' = 60")  
  
hist(recidivism$age,main='',xlab='Age at Release (Label Encoding)')
```

## 2.3. Polynomial Contrasts

## 3. Processing Day and Time Data

## 4. Processing Continuous Variables

### 4.1 Centering and Scaling

### 4.2 Transformations

### 4.3 Basis Expansions and Splines

### 4.4 Linear Projections

## Handling Missing Data

## Data Leakage

## Processing Text Data