

Nested Dummy Encoding

Andrija Djurovic

www.linkedin.com/in/andrija-djurovic

Nested Dummy VS One-Hot Encoding

- Nested dummy encoding, an alternative to standard one-hot encoding, but less commonly used in credit risk modeling.
- Nested dummy encoding applies to categorical variables that exhibit ordinal relationships among their categories.
- Nested dummy encoding creates links between adjacent categories, unlike one-hot encoding, even though in a bivariate setup both yield identical model outputs.

Nested Dummy Encoding in Credit Risk Modeling

- Establishing connections between adjacent categories makes this encoding method particularly attractive in credit risk modeling, especially when practitioners choose to bin numeric risk factors and seek statistically significant risk profiles between adjacent categories.
- This method provides opportunities to assess diverse binning methods in standard bivariate and multivariate analyses.

Constructing Nested Dummies

Assign a value of 0 to all categories below the chosen one and 1 to the selected category and those positioned above it.

Example

```
categories = c("A", "A", "B", "B", "C", "C", "C", "D", "D", "D")
```

```
nested_dummies
```

##	category_A_vs_BCD	category_AB_vs_CD	category_ABC_vs_D
## 1	0	0	0
## 2	0	0	0
## 3	1	0	0
## 4	1	0	0
## 5	1	1	0
## 6	1	1	0
## 7	1	1	0
## 8	1	1	1
## 9	1	1	1
## 10	1	1	1

Interpreting Nested Dummies

OLS regression with nested dummies

```
##  
## Call:  
## lm(formula = y1 ~ category_A_vs_BCD + category_AB_vs_CD + category_ABC_vs_D,  
## data = db)  
##  
## Coefficients:  
## (Intercept) category_A_vs_BCD category_AB_vs_CD category_ABC_vs_D  
## 0.01862 0.06086 0.05233 -0.25819
```

Target average per categorical variable

```
## A B C D  
## 0.01861955 0.07948448 0.13181361 -0.12638054
```

Coefficient replicates

```
"(Intercept)" = average["A"]  
"category_A_vs_BCD" = average["B"] - average["A"]  
"category_AB_vs_CD" = average["C"] - average["B"]  
"category_ABC_vs_D" = average["D"] - average["C"]  
  
## (Intercept) category_A_vs_BCD category_AB_vs_CD category_ABC_vs_D  
## 0.01861955 0.06086492 0.05232913 -0.25819415
```

Interpreting Nested Dummies cont.

Binomial logistic regression with nested dummies

```
##
## Call:  glm(formula = y2 ~ category_A_vs_BCD + category_AB_vs_CD + category_ABC_vs_D,
##          family = "binomial", data = db)
##
## Coefficients:
##      (Intercept)  category_A_vs_BCD  category_AB_vs_CD  category_ABC_vs_D
##           -1.5667          -0.3259           0.1926          -0.4247
##
## Degrees of Freedom: 999 Total (i.e. Null);  996 Residual
## Null Deviance:      813.5
## Residual Deviance: 808.4    AIC: 816.4
```

Log-odds of target average per categorical variable

```
##           A           B           C           D
## -1.566676 -1.892564 -1.699952 -2.124698
```

Coefficient replicates

```
"(Intercept)" = lo(average["A"])
"category_A_vs_BCD" = lo(average["B"]) - lo(average["A"])
"category_AB_vs_CD" = lo(average["C"]) - lo(average["B"])
"category_ABC_vs_D" = lo(average["D"]) - lo(average["C"])

##      (Intercept)  category_A_vs_BCD  category_AB_vs_CD  category_ABC_vs_D
##      -1.5666761    -0.3258881        0.1926122        -0.4247462
```