

1 model development

GLM model with Poisson Distribution & Hurdle Model with Binomial Distribution

Ing. Peter Tomko, M.A.

1/9/2020

```
con <- DBI::dbConnect(odbc::odbc(), "betting_ds", bigint = "integer")

select * from t_match_stats;
```

Model Estimation

Two approaches are compared in this place, namely GLM model with Poisson Distribution and hurdle model. The latter one is set with logit link function and binomial distribution - i.e. firstly the logit is estimated to predict if the the team scores and if the team scores the second part consists of Poisson distribution. It means that the special distribution is concerned when the team does not score.

```
model_vars <- c(unique(var_importance$data$Variable), "is_home", "n_goals")
master_data <- master_data %>%
  group_by(data_type) %>%
  mutate(glm_data =
    map(binned_data,
      function(i_data){
        i_data %>%
          select(one_of(model_vars)) %>%
          mutate(is_home = ifelse(is_home == 1, "yes", "no")) %>%
          rename_all(~stringr::str_replace_all(., "__", "_")) %>%
          rename_all(~stringr::str_replace_all(., "woe.", "")) %>%
          rename_all(~stringr::str_replace_all(., ".binned", "")) %>%
          as.data.frame()
      })
  )

# - calculate full model
glm_model <-
  glm(n_goals ~ .,
    family = "poisson",
    data = master_data %>%
      filter(data_type %in% "Train") %>%
      select(data_type, glm_data) %>%
      unnest(c(glm_data)) %>%
      as.data.frame() %>%
      select(-data_type))

# - select only subset of variables based on VIF
vif_glm_model <- data.frame("vif" = car::vif(glm_model))
final_vars <- rownames(vif_glm_model %>% filter(vif <= 5))
model_data <-
```

```

master_data %>%
  filter(data_type %in% "Train") %>%
  select(data_type, glm_data) %>%
  unnest(c(glm_data)) %>%
  as.data.frame() %>%
  select(-data_type) %>%
  select(n_goals, one_of(final_vars))

# - estimate final models
glm_m_final <- glm(n_goals ~ ., family = "poisson", data = model_data)

hurdle_model <- hurdle(n_goals ~ ., dist = "poisson",
  link = "logit", zero.dist = "binomial",
  data = model_data)

```

Summary of GLM Model

```

##
## Call:
## glm(formula = n_goals ~ ., family = "poisson", data = model_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2280  -1.4171  -0.2296   0.5713   4.7571
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.773e-01  3.762e-03  47.133  < 2e-16 ***
## team          -2.746e-03  7.761e-05 -35.384  < 2e-16 ***
## avg_total_goals_last_20 -2.113e-03  1.343e-04 -15.730  < 2e-16 ***
## r_ah_advantage_last_20 -1.256e-03  2.752e-04  -4.562  5.07e-06 ***
## r_ah_advantage_last_30 -9.775e-04  2.800e-04  -3.491  0.000481 ***
## league         8.681e-05  1.963e-04   0.442  0.658297
## r_team_odds_last_40    -2.620e-03  7.589e-05 -34.529  < 2e-16 ***
## r_draw_odds_last_10    -1.009e-03  1.872e-04  -5.391  7.01e-08 ***
## is_homeyes         2.186e-01  5.100e-03  42.861  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 156672  on 126247  degrees of freedom
## Residual deviance: 149071  on 126239  degrees of freedom
## AIC: 369963
##
## Number of Fisher Scoring iterations: 5

```

Summary of Hurdle Model

```

##
## Call:
## hurdle(formula = n_goals ~ ., data = model_data, dist = "poisson", zero.dist = "binomial",
##       link = "logit")
##

```

```

## Pearson residuals:
##      Min      1Q  Median      3Q      Max
## -1.5513 -0.9943 -0.2260  0.6134  7.3465
##
## Count model coefficients (truncated poisson with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      2.037e-01  5.379e-03  37.875 < 2e-16 ***
## team             -2.885e-03  1.066e-04 -27.063 < 2e-16 ***
## avg_total_goals_last_20 -2.226e-03  1.815e-04 -12.266 < 2e-16 ***
## r_ah_advantage_last_20 -1.242e-03  4.186e-04  -2.968  0.00300 **
## r_ah_advantage_last_30 -1.022e-03  4.247e-04  -2.406  0.01613 *
## league           -7.415e-04  2.694e-04  -2.752  0.00592 **
## r_team_odds_last_40    -2.502e-03  9.071e-05 -27.585 < 2e-16 ***
## r_draw_odds_last_10    -1.244e-03  2.662e-04  -4.672  2.99e-06 ***
## is_homeyes          1.947e-01  7.056e-03  27.594 < 2e-16 ***
## Zero hurdle model coefficients (binomial with logit link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.8036824  0.0090678  88.631 < 2e-16 ***
## team             -0.0045110  0.0002049 -22.011 < 2e-16 ***
## avg_total_goals_last_20 -0.0034849  0.0003571  -9.760 < 2e-16 ***
## r_ah_advantage_last_20 -0.0018979  0.0005883  -3.226  0.001254 **
## r_ah_advantage_last_30 -0.0014775  0.0006001  -2.462  0.013812 *
## league           0.0018726  0.0005062   3.699  0.000216 ***
## r_team_odds_last_40    -0.0060118  0.0003360 -17.895 < 2e-16 ***
## r_draw_odds_last_10    -0.0014270  0.0004574  -3.120  0.001810 **
## is_homeyes          0.4295940  0.0131626  32.637 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of iterations in BFGS optimization: 13
## Log-likelihood: -1.849e+05 on 18 Df

```