

Model building

Feature engineering
Prof. Dr. Jan Kirenz

The following content is based on Mine Çetinkaya-Rundel's excellent book Data Science in a Box

Feature engineering

Feature engineering

- We prefer simple models when possible, but **parsimony** does not mean sacrificing accuracy (or predictive performance) in the interest of simplicity
- Variables that go into the model and how they are represented are just as critical to success of the model
- **Feature engineering** allows us to get creative with our predictors in an effort to make them more useful for our model (to increase its predictive performance)

Same training and testing sets as before

```
# Fix random numbers by setting the seed  
# Enables analysis to be reproducible when random numbers are used  
set.seed(1116)  
  
# Put 80% of the data into the training set  
email_split <- initial_split(email, prop = 0.80)  
  
# Create data frames for the two sets:  
train_data <- training(email_split)  
test_data  <- testing(email_split)
```

A simple approach: `mutate()`

```
train_data %>%  
  mutate(  
    date = date(time),  
    dow  = wday(time),  
    month = month(time)  
  ) %>%  
  select(time, date, dow, month) %>%  
  sample_n(size = 5) # shuffle to show a variety
```

```
## # A tibble: 5 x 4  
##   time                date      dow month  
##   <dtm>              <date>    <dbl> <dbl>  
## 1 2012-03-18 17:47:58 2012-03-18      1      3  
## 2 2012-02-08 05:50:43 2012-02-08      4      2  
## 3 2012-03-16 16:06:21 2012-03-16      6      3  
## 4 2012-02-26 06:37:38 2012-02-26      1      2  
## 5 2012-01-21 16:08:13 2012-01-21      7      1
```

Modeling workflow, revisited

- Create a **recipe** for feature engineering steps to be applied to the training data
- Fit the model to the training data after these steps have been applied
- Using the model estimates from the training data, predict outcomes for the test data
- Evaluate the performance of the model on the test data

Building recipes

Initiate a recipe

```
email_rec <- recipe(
  spam ~ .,          # formula
  data = train_data  # data to use for cataloguing names and types of variables
)

summary(email_rec)
```

```
## # A tibble: 21 x 4
##   variable    type    role    source
##   <chr>      <chr>  <chr>   <chr>
## 1 to_multiple numeric predictor original
## 2 from       numeric predictor original
## 3 cc         numeric predictor original
## 4 sent_email numeric predictor original
## 5 time       date    predictor original
## 6 image      numeric predictor original
## 7 attach     numeric predictor original
## 8 dollar     numeric predictor original
## 9 winner     nominal predictor original
## 10 inherit   numeric predictor original
## 11 viagra    numeric predictor original
## 12 password  numeric predictor original
## 13 num_char  numeric predictor original
## 14 line_breaks numeric predictor original
## 15 format    numeric predictor original
## 16 re_subj    numeric predictor original
## 17 exclaim_subj numeric predictor original
## 18 urgent_subj numeric predictor original
## 19 exclaim_mess numeric predictor original
## 20 number    nominal predictor original
## 21 spam      nominal outcome  original
```


Remove certain variables

```
email_rec <- email_rec %>%  
  step_rm(from, sent_email)
```

```
## Data Recipe  
##  
## Inputs:  
##  
##   role #variables  
## outcome      1  
## predictor    20  
##  
## Operations:  
##  
## Delete terms from, sent_email
```

Feature engineer date

```
email_rec <- email_rec %>%  
  step_date(time, features = c("dow", "month")) %>%  
  step_rm(time)
```

```
## Data Recipe  
##  
## Inputs:  
##  
##   role #variables  
## outcome      1  
## predictor    20  
##  
## Operations:  
##  
## Delete terms from, sent_email  
## Date features from time  
## Delete terms time
```

Discretize numeric variables

```
email_rec <- email_rec %>%  
  step_cut(cc, attach, dollar, breaks = c(0, 1)) %>%  
  step_cut(inherit, password, breaks = c(0, 1, 5, 10, 20))
```

```
## Data Recipe  
##  
## Inputs:  
##  
##   role #variables  
## outcome      1  
## predictor     20  
##  
## Operations:  
##  
## Delete terms from, sent_email  
## Date features from time  
## Delete terms time  
## Cut numeric for cc, attach, dollar  
## Cut numeric for inherit, password
```

Create dummy variables

```
email_rec <- email_rec %>%  
  step_dummy(all_nominal(), -all_outcomes())
```

```
## Data Recipe  
##  
## Inputs:  
##  
##      role #variables  
## outcome      1  
## predictor    20  
##  
## Operations:  
##  
## Delete terms from, sent_email  
## Date features from time  
## Delete terms time  
## Cut numeric for cc, attach, dollar  
## Cut numeric for inherit, password  
## Dummy variables from all_nominal(), -all_outcomes()
```

Remove zero variance variables

Variables that contain only a single value

```
email_rec <- email_rec %>%  
  step_zv(all_predictors())
```

```
## Data Recipe  
##  
## Inputs:  
##  
##      role #variables  
## outcome      1  
## predictor    20  
##  
## Operations:  
##  
## Delete terms from, sent_email  
## Date features from time  
## Delete terms time  
## Cut numeric for cc, attach, dollar  
## Cut numeric for inherit, password  
## Dummy variables from all_nominal(), -all_outcomes()  
## Zero variance filter on all_predictors()
```

All in one place

```
email_rec <- recipe(spam ~ ., data = email) %>%  
  step_rm(from, sent_email) %>%  
  step_date(time, features = c("dow", "month")) %>%  
  step_rm(time) %>%  
  step_cut(cc, attach, dollar, breaks = c(0, 1)) %>%  
  step_cut(inherit, password, breaks = c(0, 1, 5, 10, 20)) %>%  
  step_dummy(all_nominal(), -all_outcomes()) %>%  
  step_zv(all_predictors())
```

Building workflows

Define model

```
email_mod <- logistic_reg() %>%  
  set_engine("glm")
```

```
email_mod
```

```
## Logistic Regression Model Specification (classification)  
##  
## Computational engine: glm
```


Define workflow

Workflows bring together models and recipes so that they can be easily applied to both the training and test data.

```
email_wflow <- workflow() %>%  
  add_model(email_mod) %>%  
  add_recipe(email_rec)
```

```
## == Workflow ==  
## Preprocessor: Recipe  
## Model: logistic_reg()  
##  
## — Preprocessor —  
## 7 Recipe Steps  
##  
## ● step_rm()  
## ● step_date()  
## ● step_rm()  
## ● step_cut()  
## ● step_cut()  
## ● step_dummy()  
## ● step_zv()  
##  
## — Model —  
## Logistic Regression Model Specification (classification)  
##  
## Computational engine: glm
```

Fit model to training data

```
email_fit <- email_wflow %>%  
  fit(data = train_data)
```

```
tidy(email_fit) %>% print(n = 31)
```

```
## # A tibble: 31 x 5
```

##	term	estimate	std.error	statistic	p.value
##	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
##	1 (Intercept)	-0.900	0.268	-3.35	7.94e- 4
##	2 to_multiple	-3.15	0.452	-6.97	3.25e-12
##	3 image	-1.02	0.665	-1.54	1.23e- 1
##	4 num_char	0.0408	0.0293	1.39	1.64e- 1
##	5 line_breaks	-0.00591	0.00163	-3.63	2.89e- 4
##	6 format	-0.788	0.161	-4.89	1.01e- 6
##	7 re_subj	-3.02	0.441	-6.86	7.00e-12
##	8 exclaim_subj	0.0856	0.268	0.319	7.50e- 1
##	9 urgent_subj	3.80	1.03	3.68	2.30e- 4
##	10 exclaim_mess	0.0112	0.00221	5.07	3.99e- 7
##	11 cc_X.1.68.	-0.0264	0.454	-0.0581	9.54e- 1
##	12 attach_X.1.21.	1.84	0.398	4.62	3.86e- 6
##	13 dollar_X.1.64.	-0.00909	0.230	-0.0395	9.69e- 1
##	14 winner_yes	2.08	0.408	5.10	3.36e- 7
##	15 inherit_X.1.5.	-10.4	1479.	-0.00703	9.94e- 1
##	16 inherit_X.5.10.	1.90	1.27	1.49	1.36e- 1
##	17 password_X.1.5.	-2.48	1.03	-2.41	1.59e- 2
##	18 password_X.5.10.	-13.3	816.	-0.0163	9.87e- 1
##	19 password_X.10.20.	-15.1	1149.	-0.0131	9.90e- 1
##	20 password_X.20.28.	-14.9	1329.	-0.0112	9.91e- 1
##	21 number_small	-0.662	0.168	-3.94	8.31e- 5
##	22 number_big	0.133	0.249	0.533	5.94e- 1
##	23 time_dow_Mo	-0.350	0.319	-1.10	2.72e- 1
##	24 time_dow_Di	0.101	0.283	0.357	7.21e- 1
##	25 time_dow_Mi	-0.258	0.284	-0.909	3.64e- 1
##	26 time_dow_Do	-0.123	0.285	-0.431	6.66e- 1
##	27 time_dow_Fr	0.131	0.278	0.473	6.36e- 1
##	28 time_dow_Sa	0.259	0.298	0.869	3.85e- 1
##	29 time_month_Feb	0.851	0.181	4.70	2.55e- 6
##	30 time_month_Mär	0.471	0.184	2.56	1.05e- 2
##	31 time_month_Apr	-14.0	990.	-0.0141	9.89e- 1

Make predictions for test data

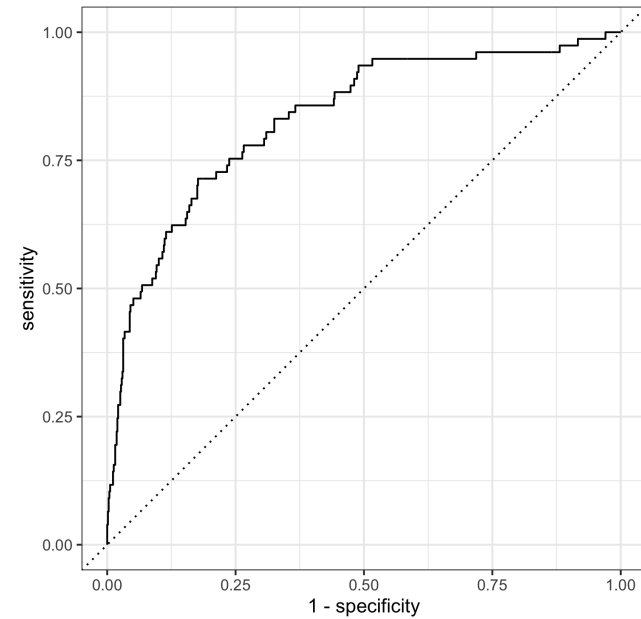
```
email_pred <- predict(email_fit, test_data, type = "prob") %>%  
  bind_cols(test_data)
```

```
email_pred
```

```
## # A tibble: 784 x 23  
##   .pred_0 .pred_1 spam  to_multiple  from    cc sent_email  
##   <dbl>   <dbl> <fct>      <dbl> <dbl> <int>    <dbl>  
## 1  0.957 0.0428  0          0      1      0        0  
## 2  0.934 0.0664  0          0      1      0        0  
## 3  0.920 0.0803  0          0      1      0        1  
## 4  0.999 0.00149 0          0      1      2        0  
## 5  0.903 0.0971  0          0      1      0        0  
## 6  0.908 0.0925  0          0      1      0        0  
## # ... with 778 more rows, and 16 more variables: time <dtm>,  
## #   image <dbl>, attach <dbl>, dollar <dbl>, winner <fct>,  
## #   inherit <dbl>, viagra <dbl>, password <dbl>, num_char <dbl>,  
## #   line_breaks <int>, format <dbl>, re_subj <dbl>,  
## #   exclaim_subj <dbl>, urgent_subj <dbl>, exclaim_mess <dbl>,  
## #   number <fct>
```

Evaluate the performance

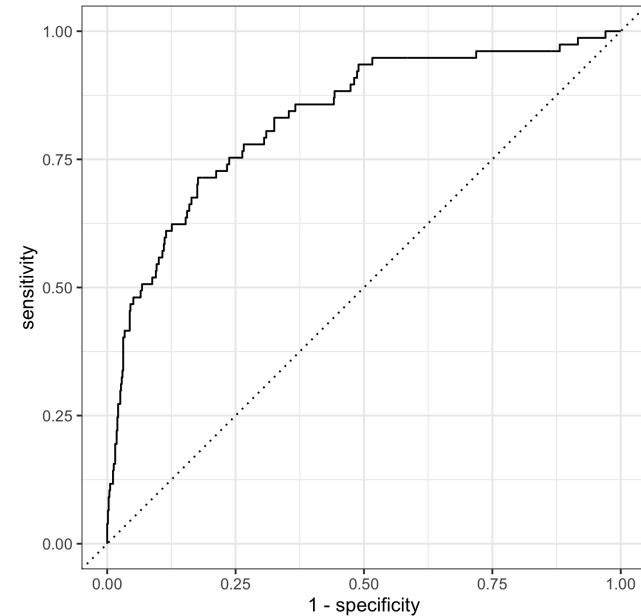
```
email_pred %>%  
  roc_curve(  
    truth = spam,  
    .pred_1,  
    event_level = "second"  
  ) %>%  
  autoplot()
```



Evaluate the performance

```
email_pred %>%  
  roc_auc(  
    truth = spam,  
    .pred_1,  
    event_level = "second"  
  )
```

```
## # A tibble: 1 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>       <dbl>  
## 1 roc_auc binary      0.831
```



Making decisions

Cutoff probability: 0.5

Output

Code

Suppose we decide to label an email as spam if the model predicts the probability of spam to be **more than 0.5**.

	Email is not spam	Email is spam
Email labelled not spam	702	68
Email labelled spam	5	9

Cutoff probability: 0.25

Output

Code

Suppose we decide to label an email as spam if the model predicts the probability of spam to be **more than 0.25**.

	Email is not spam	Email is spam
Email labelled not spam	662	40
Email labelled spam	45	37

Cutoff probability: 0.75

Output

Code

Suppose we decide to label an email as spam if the model predicts the probability of spam to be **more than 0.75**.

	Email is not spam	Email is spam
Email labelled not spam	706	72
Email labelled spam	1	5