# ETC3250/5250: Model assessment

Semester 1, 2020

Professor Di Cook

Econometrics and Business Statistics
Monash University

Week 10 (a)

```r
library(statquotes)
search_quotes(search="Holdane", fuzzy=TRUE)
```

```
## In scientific thought we adopt the simplest theory which will explain
## all the facts under consideration and enable us to predict new facts of
##     the same kind. The catch in this criterion lies in the world
## ``simplest.'' It is really an aesthetic canon such as we find implicit
## in our criticisms of poetry or painting. The layman finds such a law as
## $dx/dt = K(d^2x/dy^2)$ much less simple than "it oozes," of which it is
## the mathematical statement. The physicist reverses this judgment, and
##  his statement is certainly the more fruitful of the two, so far as
##  prediction is concerned. It is, however, a statement about something
## very unfamiliar to the plainman, namely, the rate of change of a rate
##                             of change.
## --- John Burdon Sanderson Haldane (1892--1964) Possible Worlds, 1927.
```

```
statquote(source="Box")
```

```
## It is the data that are real (they actually happened!) The model is a
   ## hypothetical conjecture that might or might not summarize and/or
              ## explain important features of the data
                       ## --- George E. P. Box
```

# Know your data. ✈️

Quantitative or qualitative response? Predictors all quantitative? Do you have independent observations?

# Know your data. ✈️

Quantitative or qualitative response? Predictors all quantitative? Do you have independent observations?

# Plot your data. 🖼️

Is there a relationship between response and predictors? Is the relationship linear? Are boundaries linear?Is variability heterogeneous? Are groups distinct? Are there unusual observations?

## Know your data. ✈️

Quantitative or qualitative response? Predictors all quantitative? Do you have independent observations?

## Plot your data. 🖼️

Is there a relationship between response and predictors? Is the relationship linear? Are boundaries linear? Is variability heterogeneous? Are groups distinct? Are there unusual observations?

## Check for missing values. 🔥

Do some variables have too many missings to use them? Do some observations have too many missings to use them? What would be a useful imputation method to fix the sporadic missing value?

# Know your data. ✈️

Quantitative or qualitative response? Predictors all quantitative? Do you have independent observations?

# Plot your data. 🖼️

Is there a relationship between response and predictors? Is the relationship linear? Are boundaries linear? Is variability heterogeneous? Are groups distinct? Are there unusual observations?
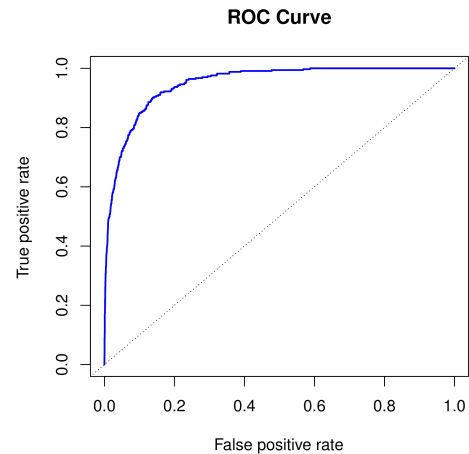
# Fit a versatile model. 💻

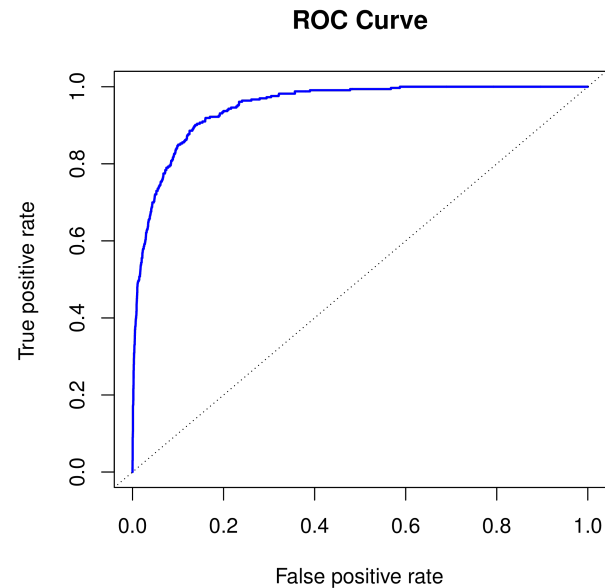Compute and plot model diagnostics. Where doesn't the model do well? How can it be refined?

# Check for missing values. 🔥

Do some variables have too many missings to use them? Do some observations have too many missings to use them? What would be a useful imputation method to fix the sporadic missing value?

# Know your data. ✈️

Quantitative or qualitative response? Predictors all quantitative? Do you have independent observations?

# Plot your data. 🖼️

Is there a relationship between response and predictors? Is the relationship linear? Are boundaries linear?Is variability heterogeneous? Are groups distinct? Are there unusual observations?

# Fit a versatile model. 💻

Compute and plot model diagnostics. Where doesn't the model do well? How can it be refined?

# Check for missing values. 🔥

Do some variables have too many missings to use them? Do some observations have too many missings to use them? What would be a useful imputation method to fix the sporadic missing value?

# ROC for classification

The ROC curve is a popular graphic for simultaneously displaying the two types of errors for all possible thresholds. It is a common method for comparing classification models. Below: ROC curve for the LDA classifier on the training set of `credit` data.
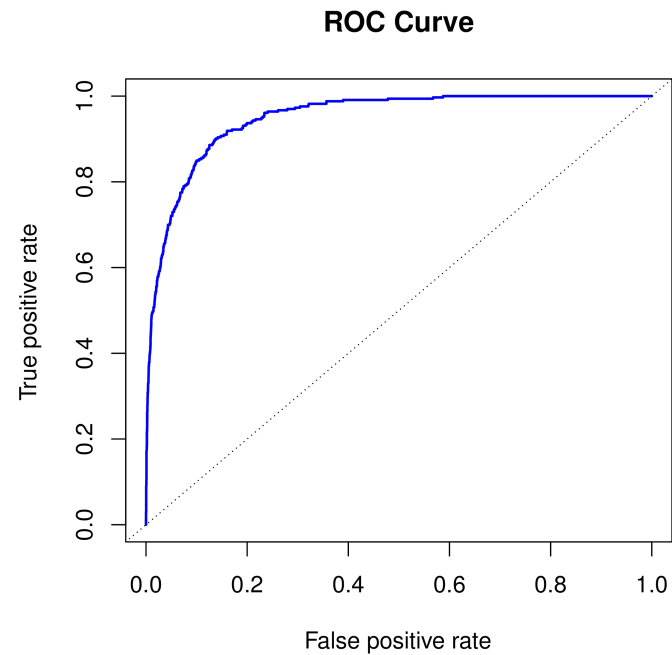
ROC Curve

## ROC Curve



The true positive rate is the sensitivity: the fraction of defaulters that are correctly identified, using a given threshold value.

The false positive rate is 1-specificity: the fraction of non-defaulters that we classify incorrectly as defaulters, using that same threshold value.

The dotted line is "no information" classifier; class and predictor are not associated.

The ideal ROC curve hugs the top left corner, indicating a high true positive rate and a low false positive rate.

## ROC Curve



If the classifier returns a prediction between 0 and 1, interpret as the probability of a positive, then threshold (split the data) at different values, e.g. 0.1, 0.2, 0.3, 0.4, 0.5, ...

Compute the confusion table for each split, record the sensitivity and specificity and plot the resulting numbers.

|  | | true | |
|---|---|---|---|
|  | | C1 (positive) | C2 (negative) |
| pred- | C1 | *a* | *b* |
| icted | C2 | *c* | *d* |

📊 Sensitivity: *a/(a+c)* (true positive, recall)
📊 Specificity: *d/(b+d)* (true negative)

```
library(tidyverse)
library(yardstick)
glimpse(two_class_example)
```

```
## Rows: 500
## Columns: 4
## $ truth     <fct> Class2, Class1, Class2, Class1
## $ Class1    <dbl> 0.0035892426, 0.6786210540, 0.
## $ Class2    <dbl> 9.964108e-01, 3.213789e-01, 8.
## $ predicted <fct> Class2, Class1, Class2, Class1
```

Set threshold to 0.5

```
##           Truth
## Prediction Class1 Class2
##     Class1    227     31
##     Class2     50    192
```

sensitivity = 0.82, 1-specificity = 0.14



Your turn: Set the threshold to be 0.75, re-compute the confusion matrix, and sensitivity, specificity.

Really nice explanation by Parul Pandey here:

# ROC for classification



(left) LDA and SVM similar.
(right) SVM radial basis with $\gamma = 10^{-1}$ is the best.

Fig 9.10

# Multivariate outliers

Mahalanobis distance measures the distance from the mean, relative to the variance-covariance matrix, and is useful for outlier detection:
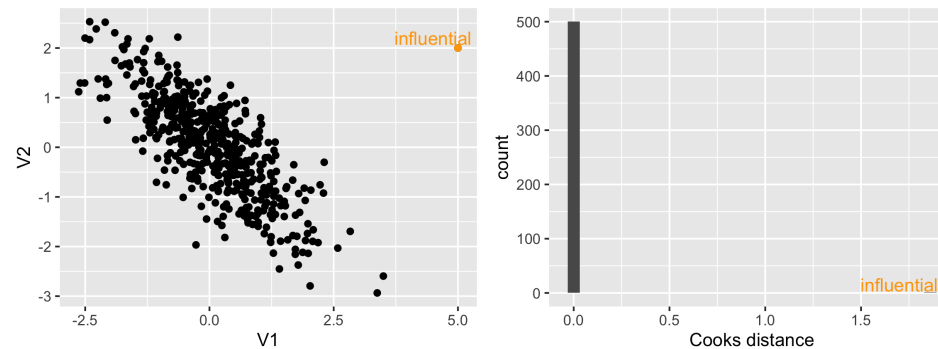$$D^2 = (X - \mu)'\Sigma^{-1}(X - \mu)$$



Related to "leverage" in regression diagnostics.

# Influential observations

Cook's distance measures the change in the model estimates due to the observation: $D_i = \frac{e_i^2}{MSE \times p} \frac{h_i}{(1-h_i)^2}$ where $h_i$ is the leverage of observation $i$.

# Utilising bagging

Remember the vote matrix available from random forests:

$$V = (V_1 V_2 \ldots V_K)$$

$$= \begin{bmatrix} p_{11} p_{12} & \cdots & p_{1K} \\ p_{21} p_{22} & \cdots & p_{2K} \\ \cdots \cdots & & \cdots \\ p_{n1} p_{n2} & \cdots & p_{nK} \end{bmatrix}$$

With bagging, multiple out of bag predictions produces uncertainty measure for each observation. It's possible that observations with higher uncertainty are outliers.

# Variable importance

📊 Working with standardised variables helps, because magnitude of coefficients is then directly interpreted as importance
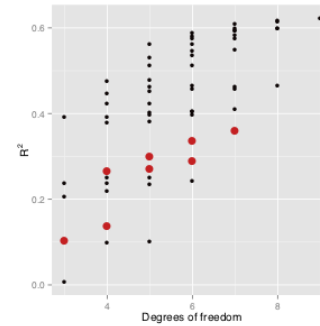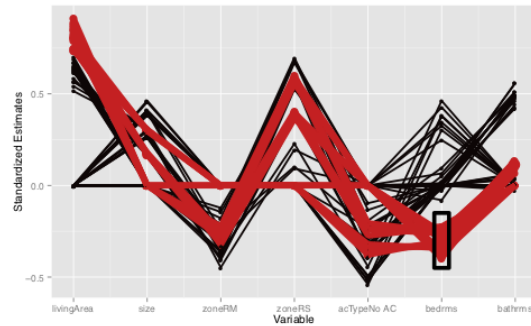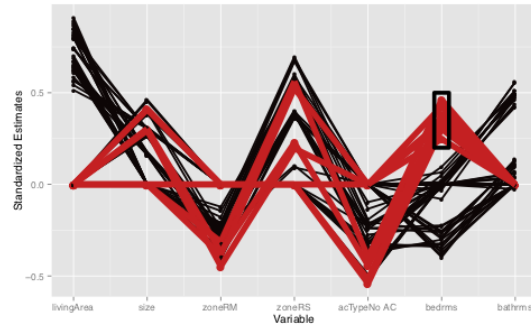
📊 Permutation approach in random forests is useful more broadly. Compare magnitude of coefficients between models built on original and permuted variable.

📊 Effect of one predictor with the response can depend on their relationship with one another. Called multicollinearity in regression.

# Bigger picture

All possible model fits to housing data with 7 variables, from Wickham et al (2015) Removing the Blindfold

Three typical estimates for bedrooms: big positive, close to 0, big negative.

Models with big positive coefficients for bedrooms tend to have weaker fits. They tend to occur with models that have no livingArea contribution, and more negative coefficients for zoneRM, and no air con.

Models with big negative coefficients on bedrooms tend to have stronger fits. All contrast with livingArea (high positive coefficients).

If bedrooms contribute to the model, bathrooms do not.

# Model choice - robustness of conclusions

Whatever way you model the data, the interpretations should be consistent.

📊 Bias can explain difference in predictions between models, flexible vs inflexible can provide a spectrum on what the data predicts.
📊 Broad changes in a model when some cases or some variables are not used, should evoke suspicions (your "spidey sense").
📊 Model fit statistics are a measure of predictive power. A weak model can still be useful if there is a large cost involved.

# 🧑‍💻 Made by a human with a computer

Slides at https://iml.numbat.space.

Code and data at https://github.com/numbats/iml.

Created using R Markdown with flair by **xaringan**, and **kunoichi** (female ninja) style.