

ETC3250/5250 Introduction to Machine Learning

Week 7: Explainable artificial intelligence (XAI)

Professor Di Cook

etc3250.clayton-x@monash.edu

Department of Econometrics and Business Statistics

Overview

We will cover:

- Global explainability
- Local explainability
 - LIME
 - Counterfactuals
 - Anchors
 - Shapley values

Global explainability

Variable importance (1/4)

Remember:

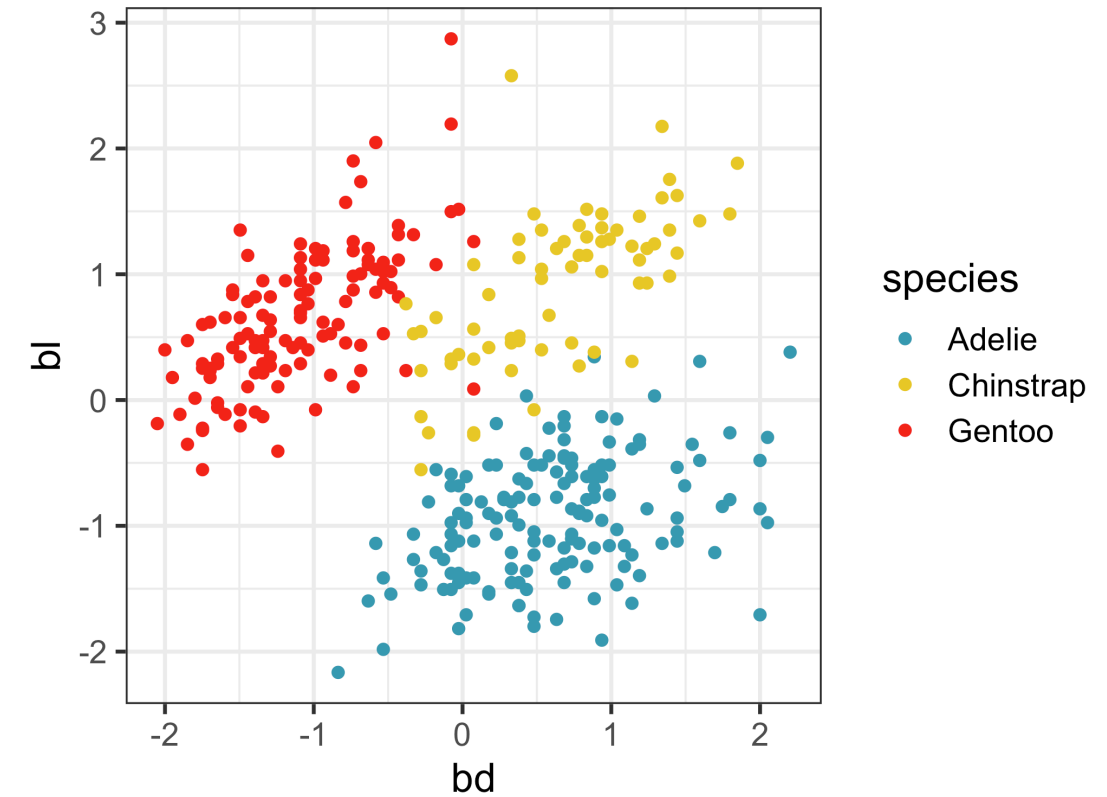
- Model coefficients on standardised data
- Effect of **collinearity**
- Importance from permutation

Variable importance (2/4)

Model coefficients on **standardised data**

```
parsnip model object  
  
Call:  
lda(species ~ ., data = data, prior = ~c(1/3, 1/3, 1/3))  
  
Prior probabilities of groups:  
  Adelie Chinstrap  Gentoo  
  0.33    0.33    0.33  
  
Group means:  
      b1    bd    fl    bm  
Adelie -0.94  0.61 -0.78 -0.62  
Chinstrap 0.90  0.64 -0.36 -0.58  
Gentoo 0.66 -1.10  1.16  1.09
```

b1 and **bd** are the most important variables

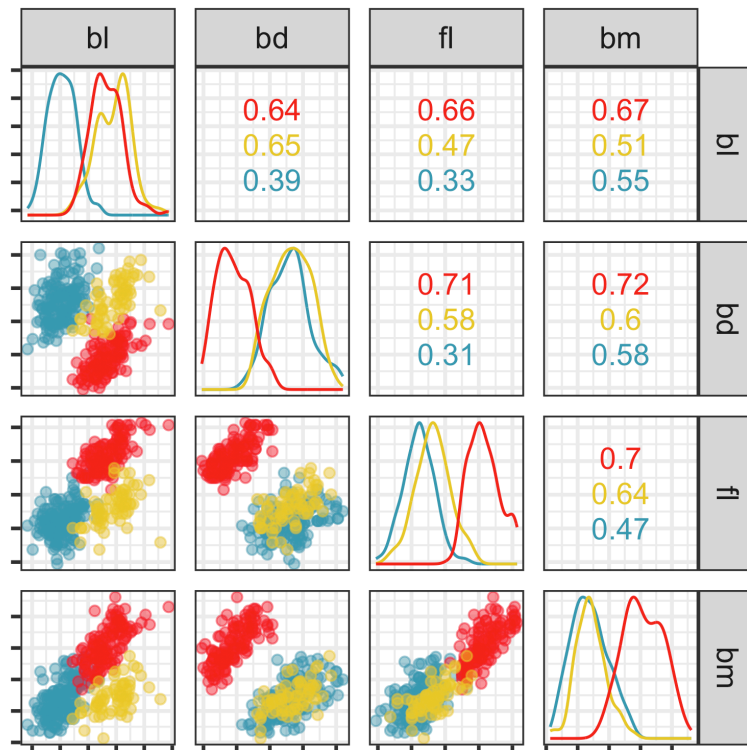


Variable importance (3/4)

When predictors are strongly linearly associated, interpreting coefficients purely on magnitude can be incorrect.

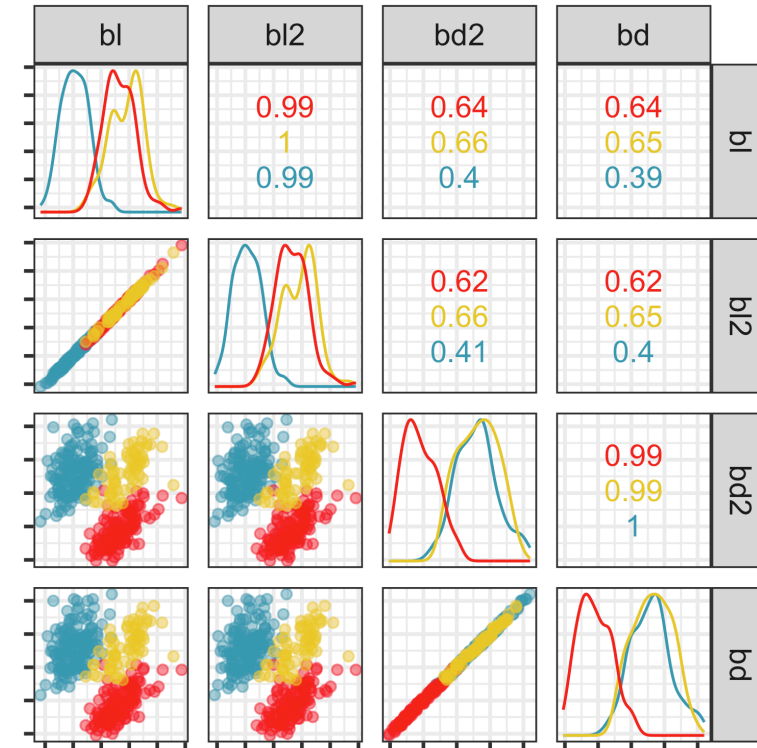
Original

	LD1	LD2
b1	-0.24	-2.319
bd	2.04	0.172
f1	-1.22	0.062
bm	-1.18	1.257



Correlated variables

	LD1	LD2
b1	1.80	-1.921
b12	-1.57	-0.407
bd2	0.49	1.548
bd	-2.54	-1.375
f1	1.21	0.052
bm	1.21	1.270



Permutation variable importance (1/2)

For trained model \hat{f} , which depends on data \mathbf{X} to predict response y , with loss function $L(y, \hat{f})$ (e.g. misclassification rate, error),

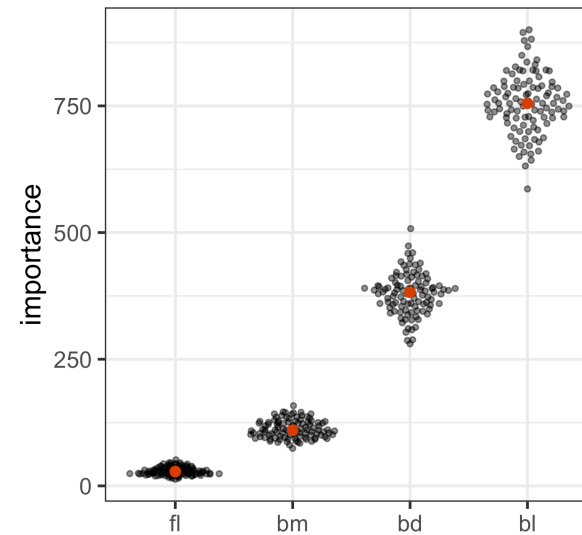
1. Estimate $L(y, \hat{f})$ on the data, L^{orig} .
2. For each variable $j \in 1, \dots, p$,
 - Generate data matrix \mathbf{X}^{perm} by **permuting** variable j . **This breaks the association** between variable j and observed y .
 - Compute the $L(y, \hat{f})$ on the permuted data, L^{perm} .
 - Compare L^{orig} and L^{perm} , e.g. $|L^{\text{orig}} - L^{\text{perm}}|$
3. Most important variables have larger values.

Permutation variable importance (2/2)

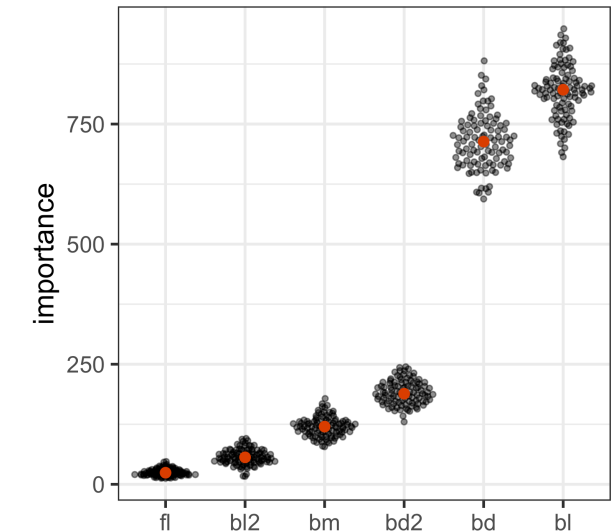
Random forests have this baked into the model fitting (using the out-of-bag cases).

Generally, should be conducted on the test set.

```
1 # Using DALEX with tidymodels
2 # https://www.tmwr.org/explain
3 # https://ema.drwhy.ai/featureImportance.html
4 vip_features <- colnames(p_std)[2:5]
5
6 vip_train <-
7   p_std |>
8   select(all_of(vip_features))
9
10 explainer_lda <-
11   explain_tidymodels(
12     lda_fit,
13     data = vip_train,
14     y = p_std$species,
15     verbose = FALSE
16   )
17 vip_lda <- model_parts(explainer_lda,
18                       B=100)
```



Data with additional correlated variables



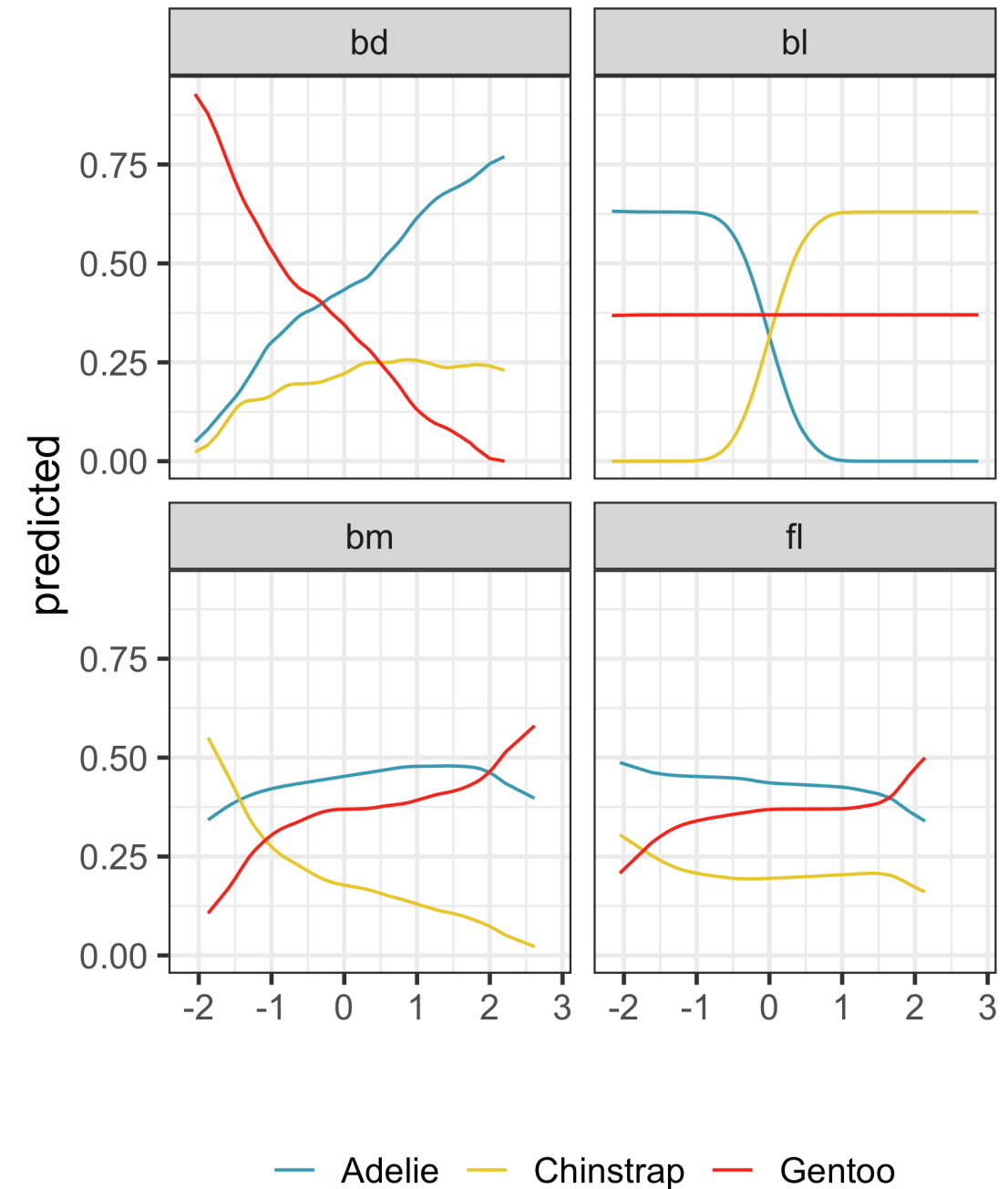
- Variables with correlation still can affect results.
- Variables can mask the importance of others.

Partial dependence profiles (1/2)

Partial dependence profiles show how the model prediction changes across different values of an explanatory variable.

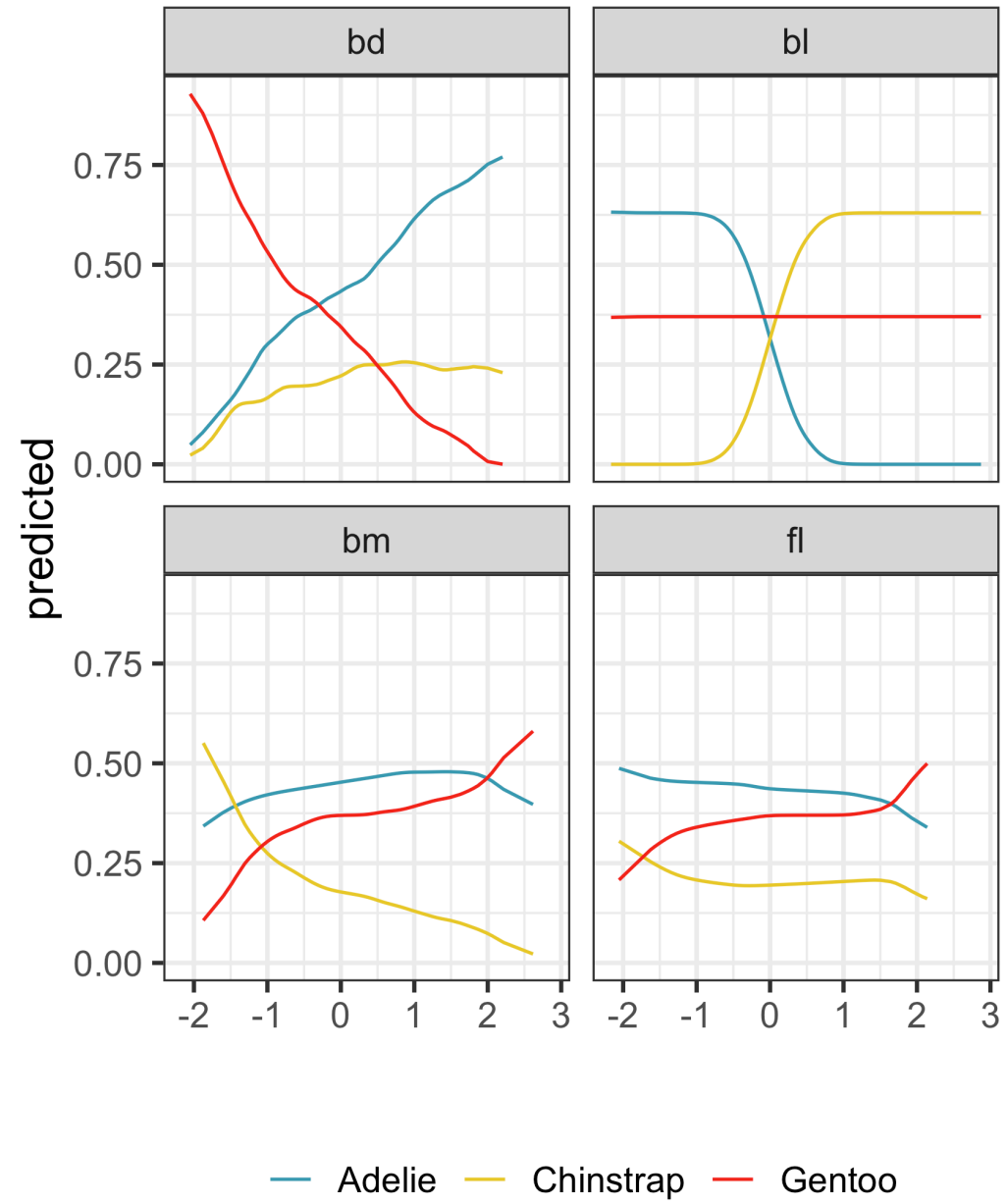
```
1 # With DALEX
2 pdp_lda <- model_profile(
3     explainer_lda,
4     N=100)
```

Shows what the model sees.

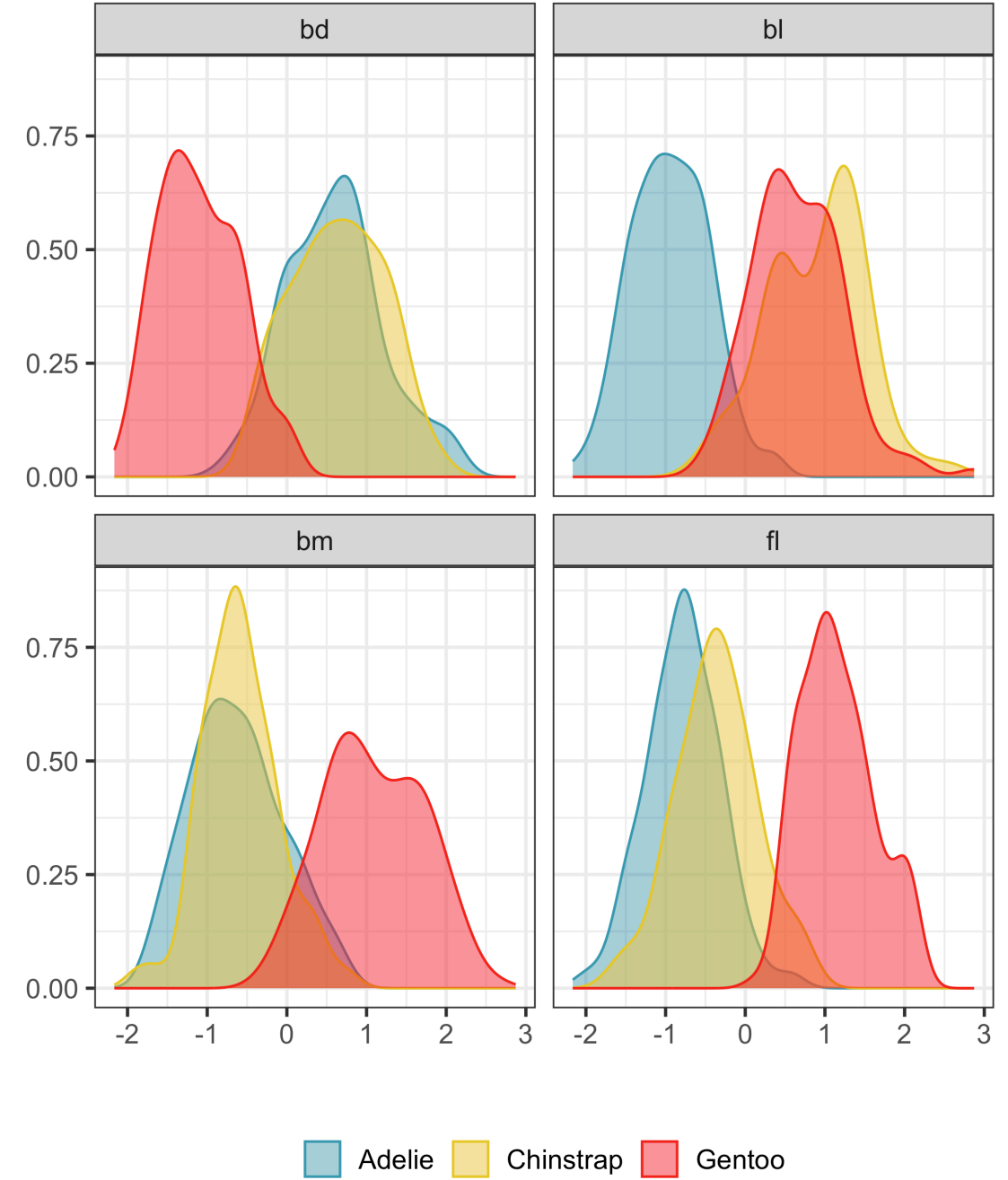


Partial dependence profiles (1/2)

PDP suggests LDA sees

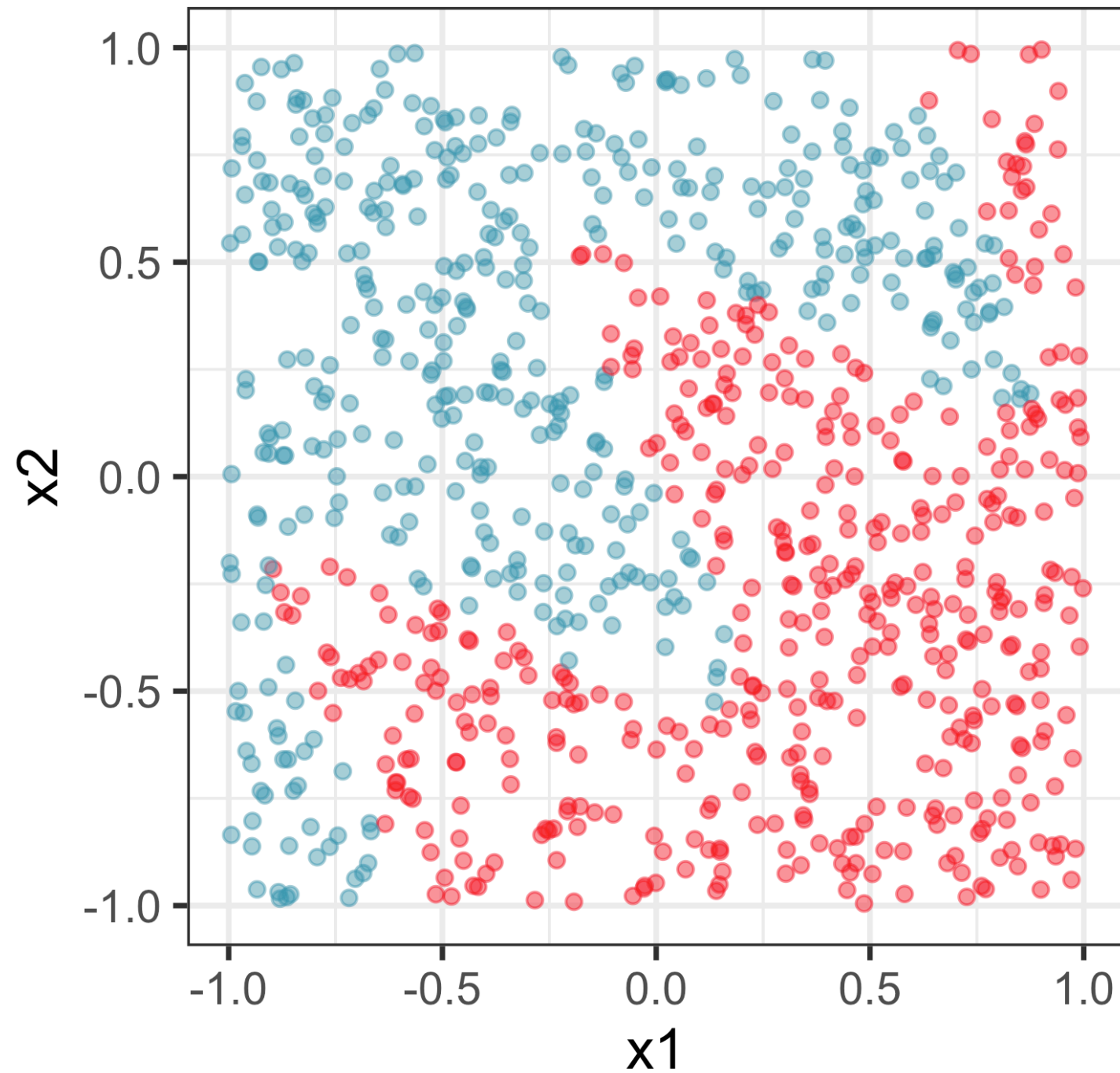


What do we see?



Local explainability

Linear vs non-linear separation

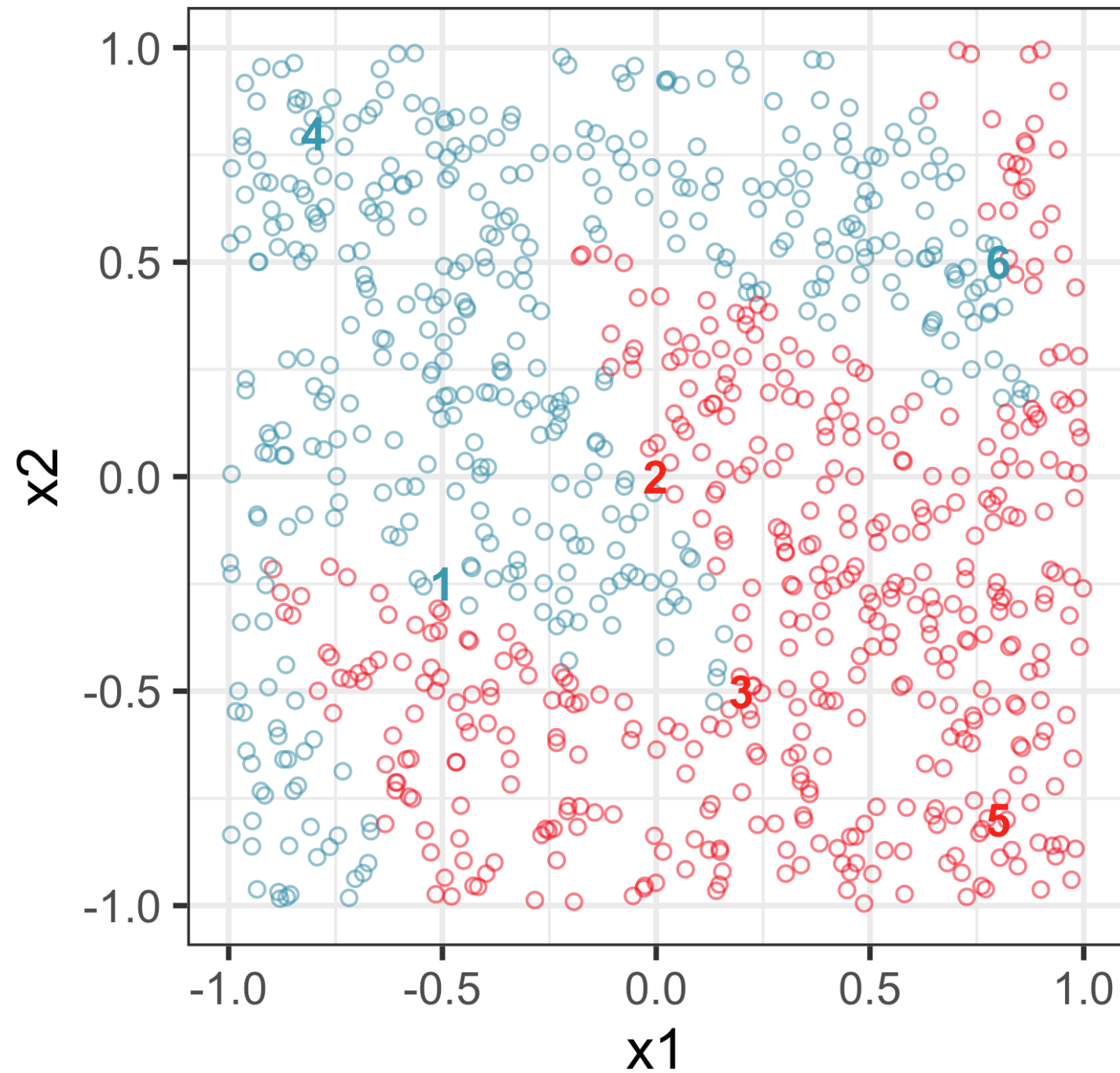


When the difference between classes is non-linear, **variable importance changes locally**.

Mark a point where x_1 is most important in distinguishing the classes.

Mark a point where x_2 is most important in distinguishing the classes.

Selected points to use for illustration



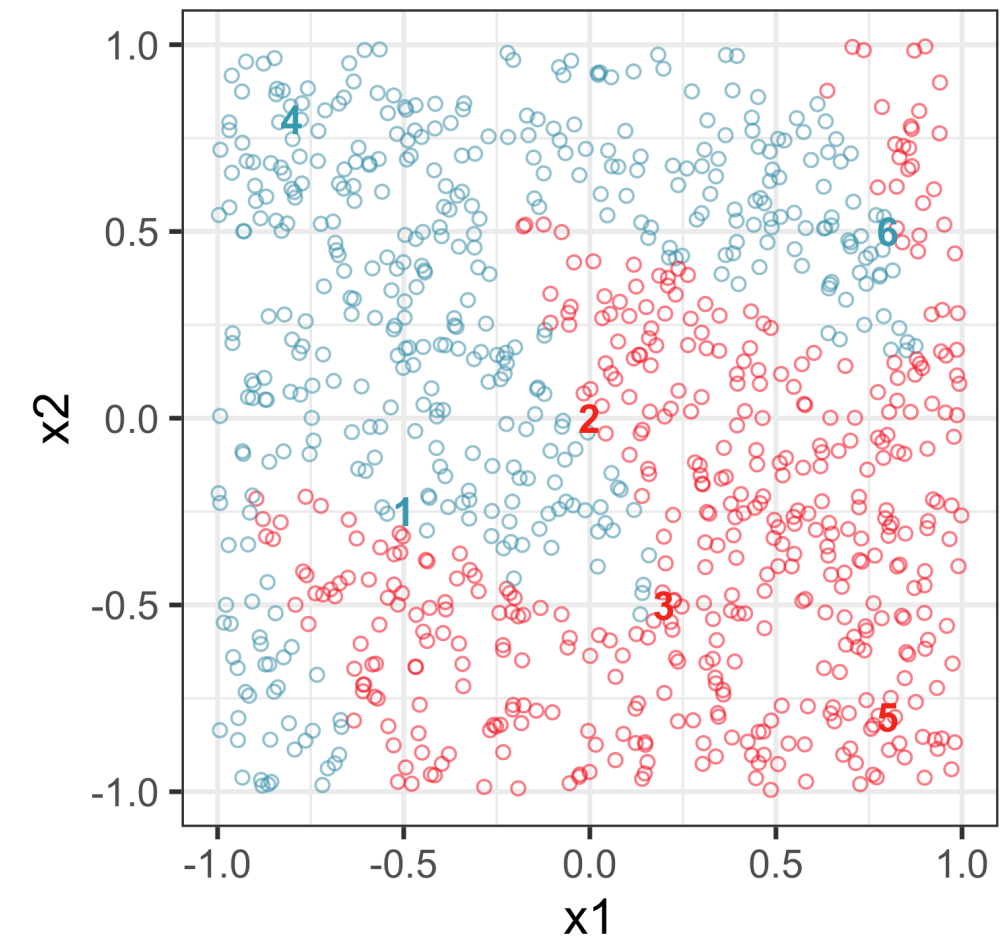
Which variable is most important?

obs	expect
1	x_1
2	x_2
3	x_2 ?
4	x_1, x_2
5	x_1, x_2
6	x_2

LIME

Fit a linear regression in the local neighbourhood of observation of interest.

```
1 library(DALEXtra)
2 library(lime)
3 w_rf <- randomForest(cl~., data=w)
4 w_rf_exp <- DALEX::explain(model = w_rf,
5                           data = w[, 1:2],
6                           y = w$cl == "A")
7 model_type.dalex_explainer <-
8   DALEXtra::model_type.dalex_explainer
9 predict_model.dalex_explainer <-
10  DALEXtra::predict_model.dalex_explainer
11 w_lime <- predict_surrogate(
12   explainer = w_rf_exp,
13   new_observation = w_new,
14   n_features = 2,
15   n_permutations = 100,
16   type = "lime")
```



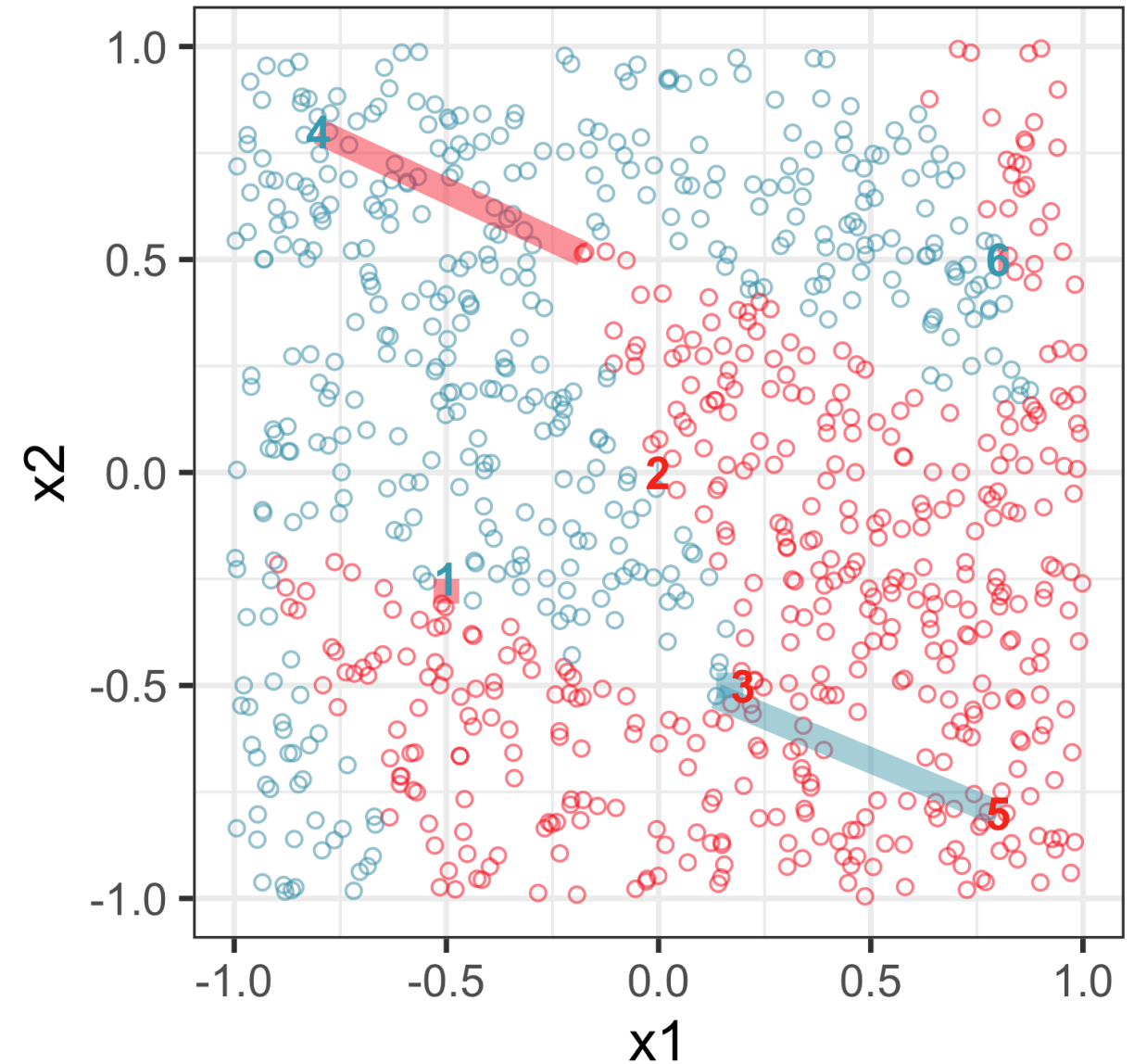
```
# A tibble: 6 × 4
  case model_intercept  x1    x2
<chr> <dbl> <dbl> <dbl>
1 1      0.485 -0.419 0.189
2 2      0.643 -0.184 -0.00800
3 3      0.474 0.141 0.347
4 4      0.695 -0.354 -0.441
5 5      0.466 0.212 0.387
6 6      0.498 0.339 -0.379
```

Counterfactuals

Find the **closest observation (counterfactual)** that has the **different class**.
What values of the variables would you need to change to change the observation of interest into the counterfactual.

```
1 library(iml)
2 # devtools::install_github("dandls/counterfactuals")
3 library(counterfactuals)
4 predictor_rf = iml::Predictor$new(w_rf,
5                                 type = "prob")
6 # predictor_rf$predict(w_new[1,])
7 w_classif <- counterfactuals::NICEClassif$new(
8   predictor_rf)
9
10 w_new_cf <- w_new
11 w_new_cf$c1 <- ifelse(w_new[,3]=="A",
12                    "B", "A")
13 for (i in 1:nrow(w_new)) {
14   w_cf = w_classif$find_counterfactuals(
15     x_interest = w_new[i,],
16     desired_class = w_new_cf[i,3],
17     desired_prob = c(0.5, 1)
18   )
19   w_new_cf[i, 1] <- w_cf$data$x1
20   w_new_cf[i, 2] <- w_cf$data$x2
21 }
```

	x1o	x2o	c1o	x1	x2	c1
1	-0.5	-0.25	A	-0.5000	-0.31	B
2	0.0	0.00	B	-0.0057	0.00	A
3	0.2	-0.50	B	0.1358	-0.50	A
4	-0.8	0.80	A	-0.1785	0.51	B
5	0.8	-0.80	B	0.1358	-0.52	A
6	0.8	0.50	A	0.8249	0.50	B

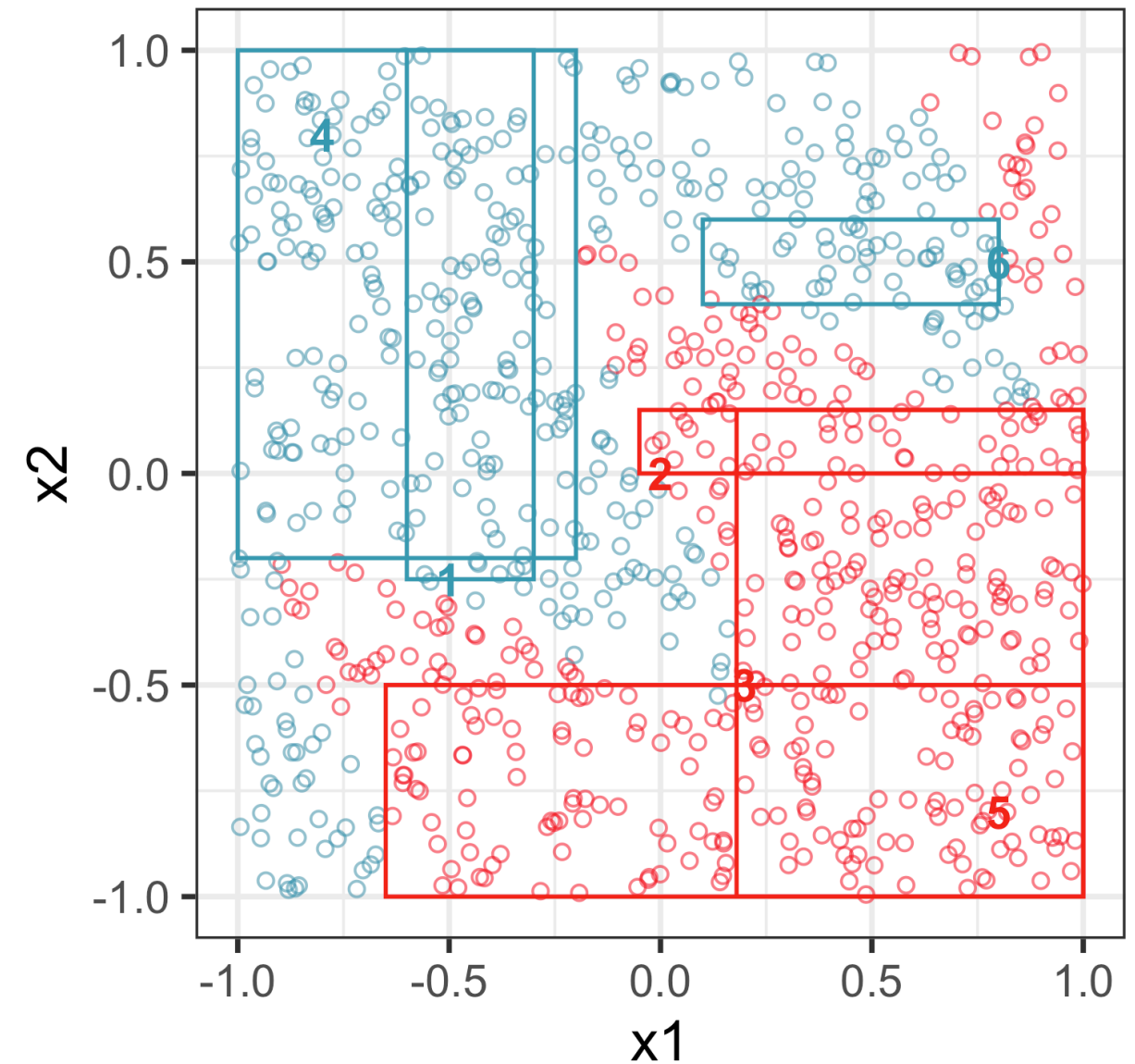


Note: If case is misclassified, the desired class needs to be the true class.

Anchors

How far can you extend from the value of the observation in each direction and still have all observations be the same class.

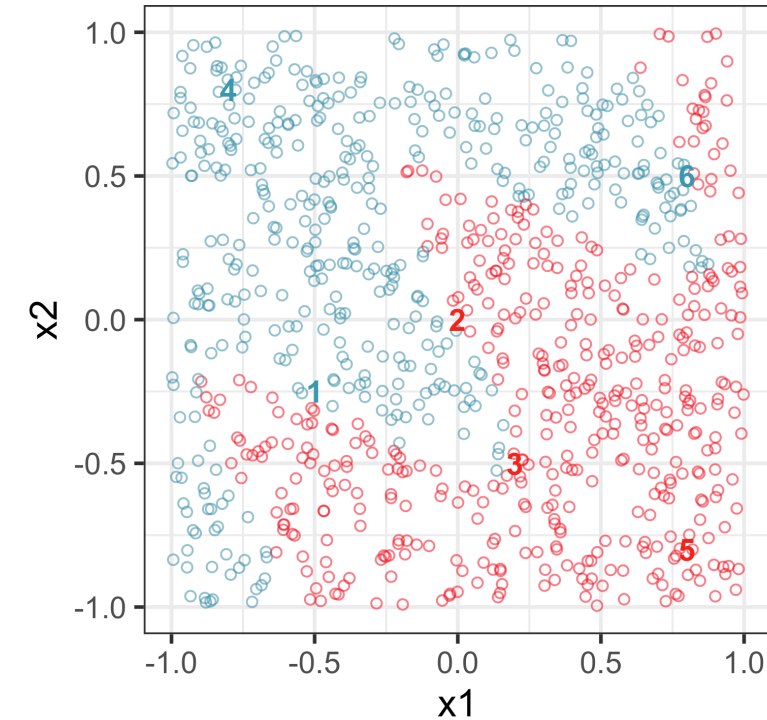
Note: No working R package to calculate these.



Shapley values

A Shapley value is computed from the change in prediction when all combinations of presence or absence of other variables. In the computation, for each combination, the prediction is computed by substituting absent variables with their average value.

```
1 library(kernelshap)
2 library(shapviz)
3 w_explain <- kernelshap(
4   w_rf,
5   w_new[,1:2],
6   w[,1:2],
7   verbose = FALSE
8 )
```



	x1	x2	cl	shapAx1	shapAx2
1	-0.5	-0.25	A	0.358	0.15
2	0.0	0.00	B	-0.236	-0.25
3	0.2	-0.50	B	-0.164	-0.32
4	-0.8	0.80	A	0.255	0.26
5	0.8	-0.80	B	-0.215	-0.27
6	0.8	0.50	A	-0.059	0.57

Summary

Which variable is most important?

obs	expect	LIME	CF	SHAP
1	x1	x1	x2	x1
2	x2	x1	x1	x1, x2
3	x2 ?	x2	x1	x2
4	x1, x2	x1, x2	x1, x2	x1, x2
5	x1, x2	x2	x1, x2	x1, x2
6	x2	x2	x1	x2

They don't all agree.

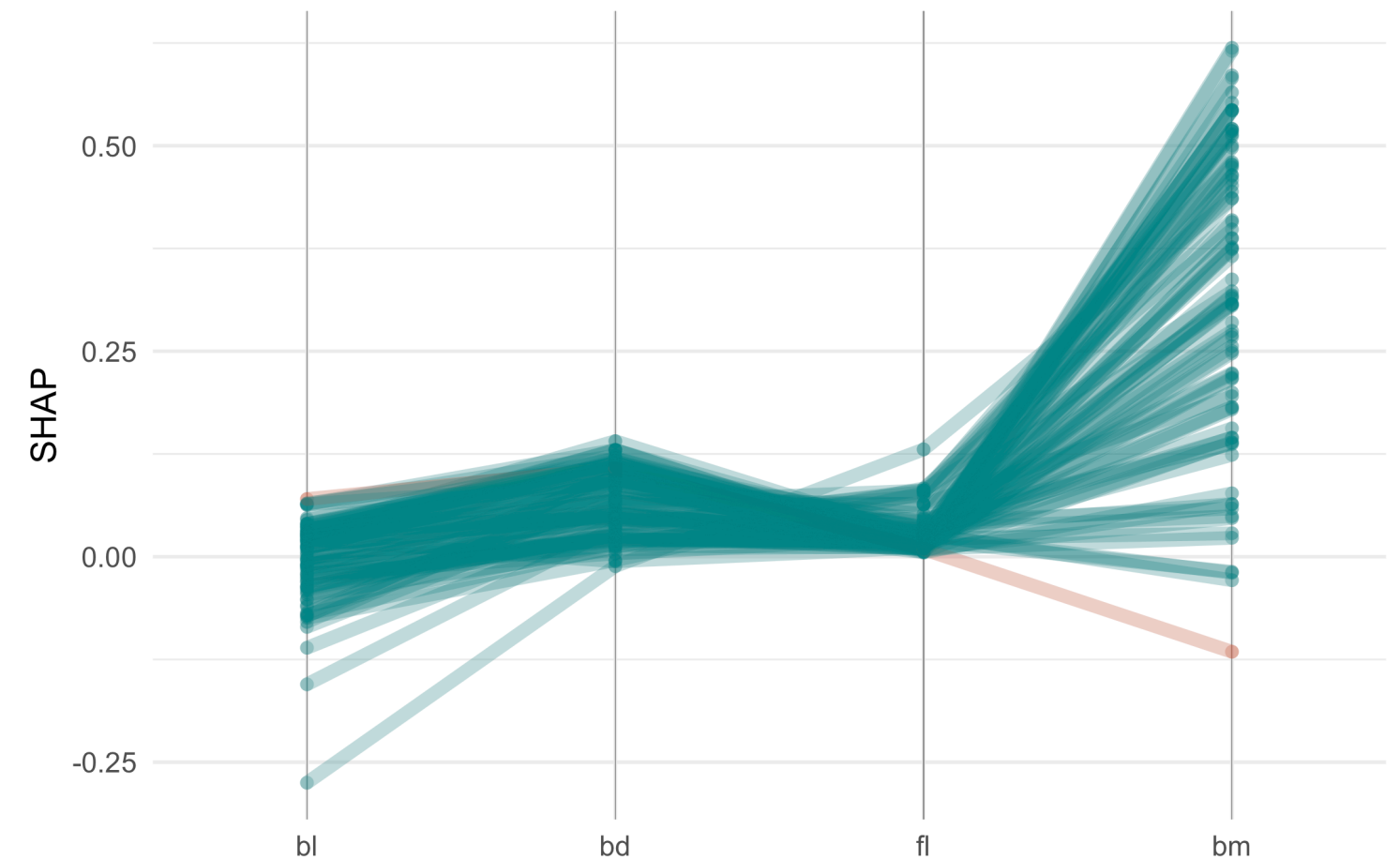
You **need good visualisation of the model in the data space** to fully digest the importance of the variables.

Example: penguins (1/2)

Compute SHAP values for the neural network model

```
1 library(keras)
2 p_nn_model <- load_model_tf("../data/penguins_
3 p_nn_model
4
5 # Explanations
6 # https://www.r-bloggers.com/2022/08/kernel-sh
7 library(kernelshap)
8 library(shapviz)
9 p_explain <- kernelshap(
10   p_nn_model,
11   p_train_x,
12   bg_X = p_train_x,
13   verbose = FALSE
14 )
15 p_exp_sv <- shapviz(p_explain)
16 save(p_exp_sv, file="../data/p_exp_sv.rda")
```

Highlight SHAP values for a misclassified Gentoo penguin



Note: the SHAP value is much lower than values for all other penguins on **bm**.

Example: penguins (2/2)

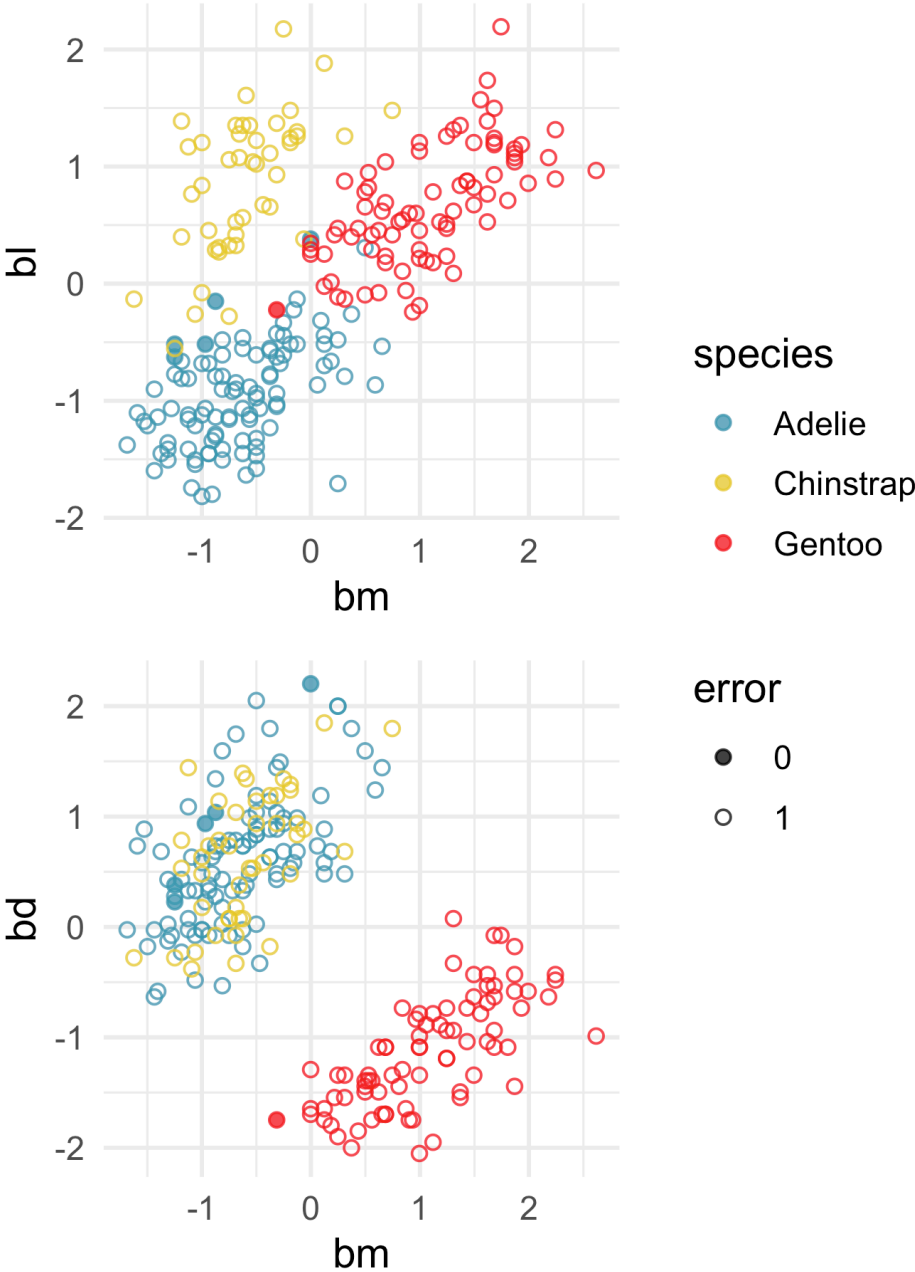
Weights from hidden layer

	[,1]	[,2]
[1,]	0.56	0.80
[2,]	0.17	-0.21
[3,]	-0.15	-0.15
[4,]	-0.80	0.54

Model uses mostly **bl** and **bm**.

Note: this analysis used the **training set** because this Gentoo penguin was misclassified as an Adelie in the training set.

	p_train_pred_cat		
	Adelie	Chinstrap	Gentoo
Adelie	95	5	0
Chinstrap	0	45	0
Gentoo	1	0	81



**Next: Support vector machines
and nearest neighbours**