

GO 
DATA
DRIVEN

Unpacking the Black Box

How to Interpret your Machine Learning Model

About us



Marysia Winkels

Data Science Educator @ GoDataDriven



James Hayward

Data Science Educator @ GoDataDriven

GoDataDriven

GO
DATA
DRIVEN



Organize

Data and AI Strategy Execution



Build

Data Platforms and AI Applications



Train

Data and AI Learning Journeys



Manage

Managed Data Platform and AI Services



Democratize

Data Democratization



Solve

Turn-Key Cloud & Data Solutions

Trusted by global industry leaders



Uber

DUPONT™

Booking.com



AIRBUS

MERCK

Ahold
Delhaize

CREDIT SUISSE



ING The ING lion logo is a detailed orange illustration of a lion standing on its hind legs.

TRADE MARK
Heineken®

BOSCH

rbi reed business
information



ABN·AMRO

TOMTOM® The TomTom logo features the brand name in a bold, black, sans-serif font next to a silver circular icon containing a red hand.

The Verizon logo consists of a red checkmark shape above the word "verizon" in a bold, black, sans-serif font.

AIR FRANCE KLM

intel®

GoDataDriven Academy

A variety of courses

- Data Science with Python
(introductory & advanced)
- Deep Learning
- Data Visualisation & Storytelling
- Analytics Translation
- Time Series
- Machine Learning Explainability

GO 
DATA
DRIVEN



About you!

Join us at PollEverywhere

GO 
DATA
DRIVEN



gdd.li/poll-amld

Workshop outline

Topics are covered in the following way:

1. Experiment *first*
2. Explain *later*

Put it in practice:

- Hackathon
- Lightning talks!



Workshop content

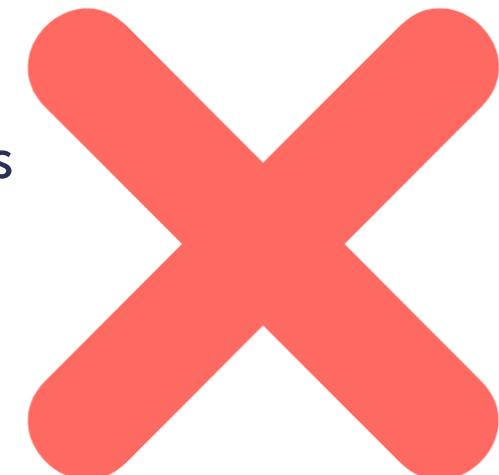
What we will cover:

- Intrinsicly explainable machine learning models
- Methods for model-agnostic interpretability
- Advantages and disadvantages of each method



What this workshop is **not** about:

- The theory behind every machine learning model (*algorithmic transparency*)
- Shapley values
- Neural networks



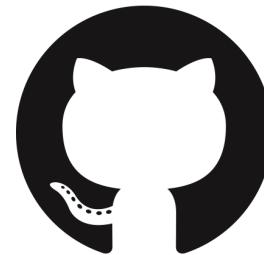
Material

Two ways to access the material:

1. Binder

*Does not require any
installation*

2. Download or clone
the material from our
GitHub repository



gdd.li/github

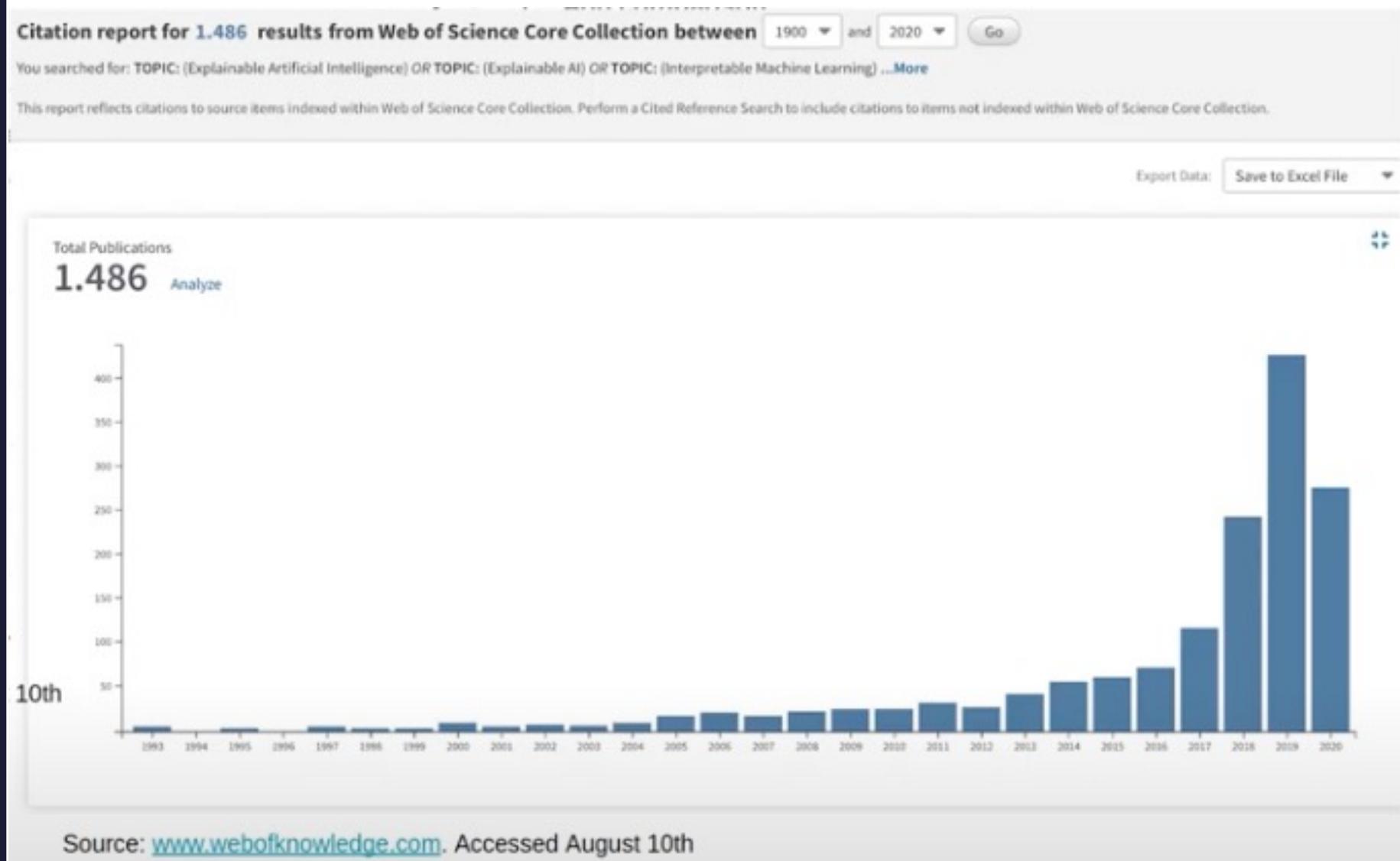


gdd.li/binder

Unpacking the black box

History

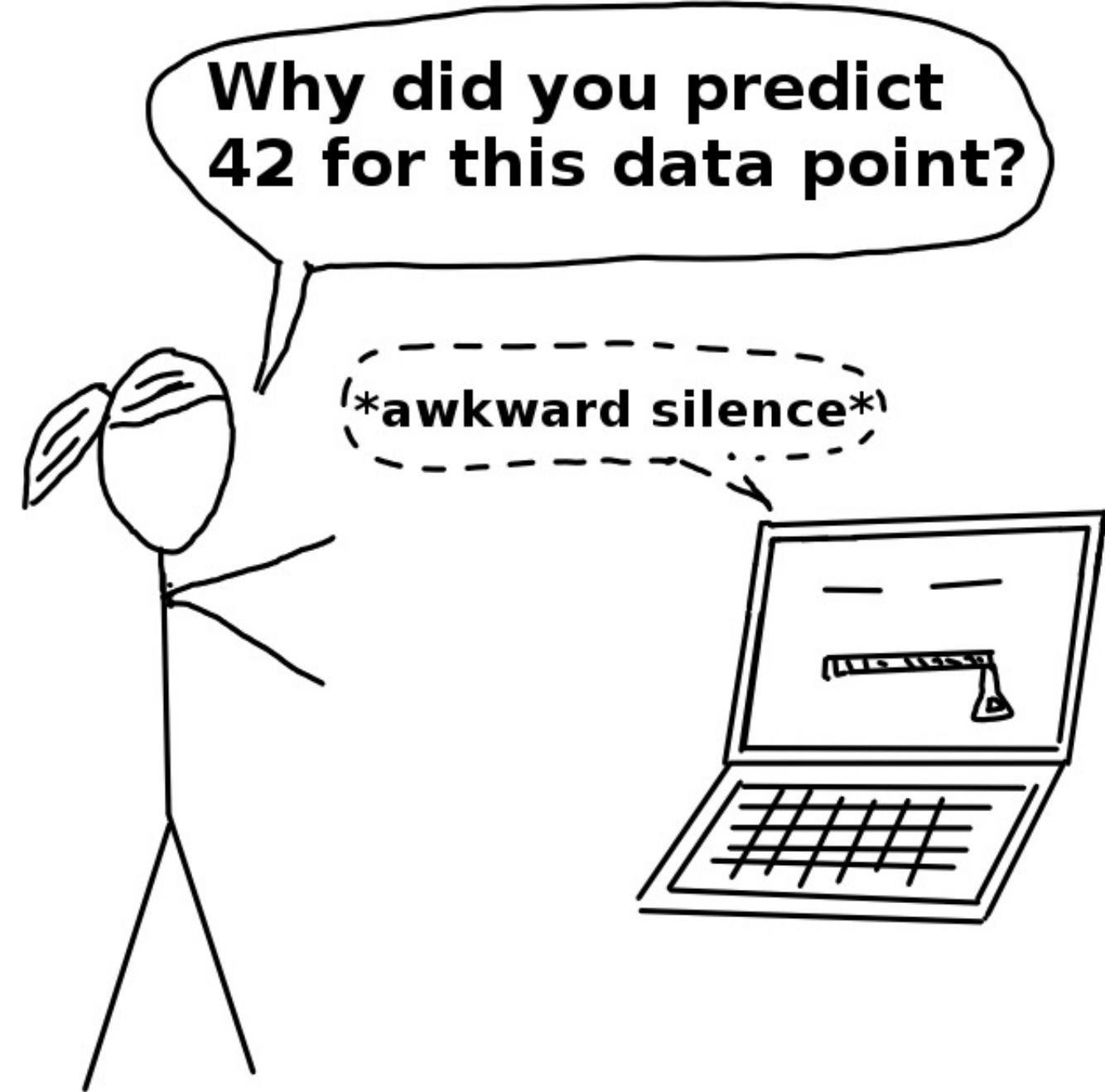
- 1800s – Linear Regression
- 1960s – Decision Tree
- 2001 – Random Forests
- 2012 – Deep Learning
- 2016 – *Interpretable Machine Learning*



“Interpretability is the degree to which a human can understand the cause of a decision”

Source: Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences.

Why of Interpretability



Why of Interpretability

- Engender trust



Why of Interpretability

- Engender trust
- Right to explanation

Source:

<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Two Petty Theft Arrests

VERNON PRATER

Prior Offenses

2 armed robberies, 1 attempted armed robbery

Subsequent Offenses

1 grand theft

LOW RISK

3

BRISHA BORDEN

Prior Offenses

4 juvenile misdemeanors

Subsequent Offenses

None

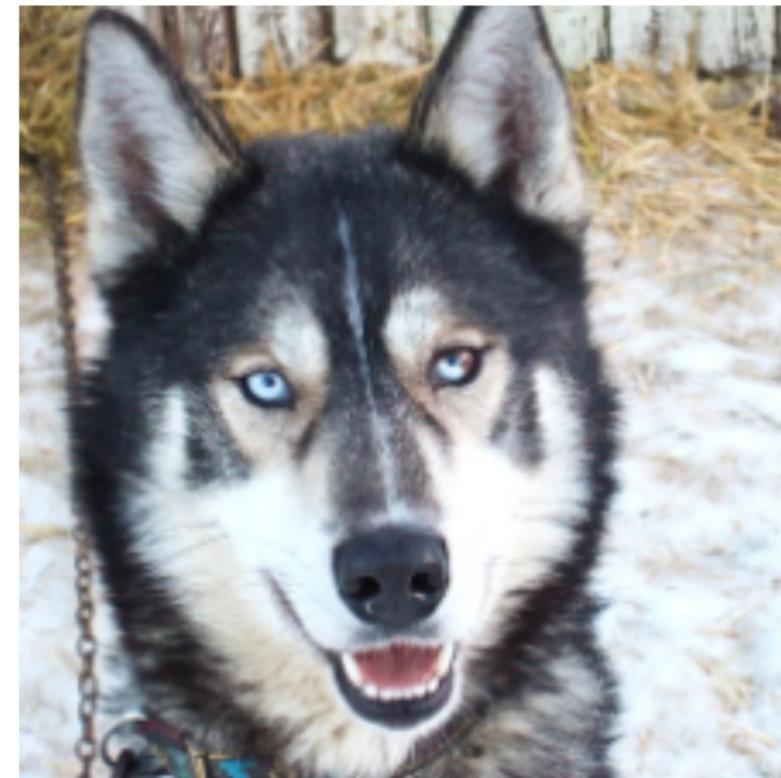
HIGH RISK

8

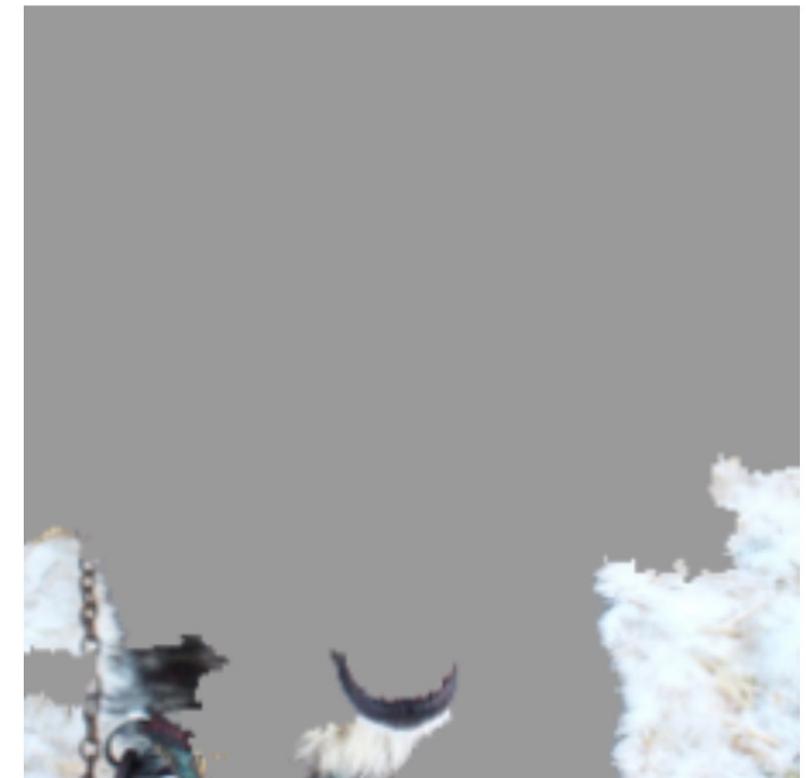
Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Why of Interpretability

- Engender trust
- Right to explanation
- Debug models



(a) Husky classified as wolf

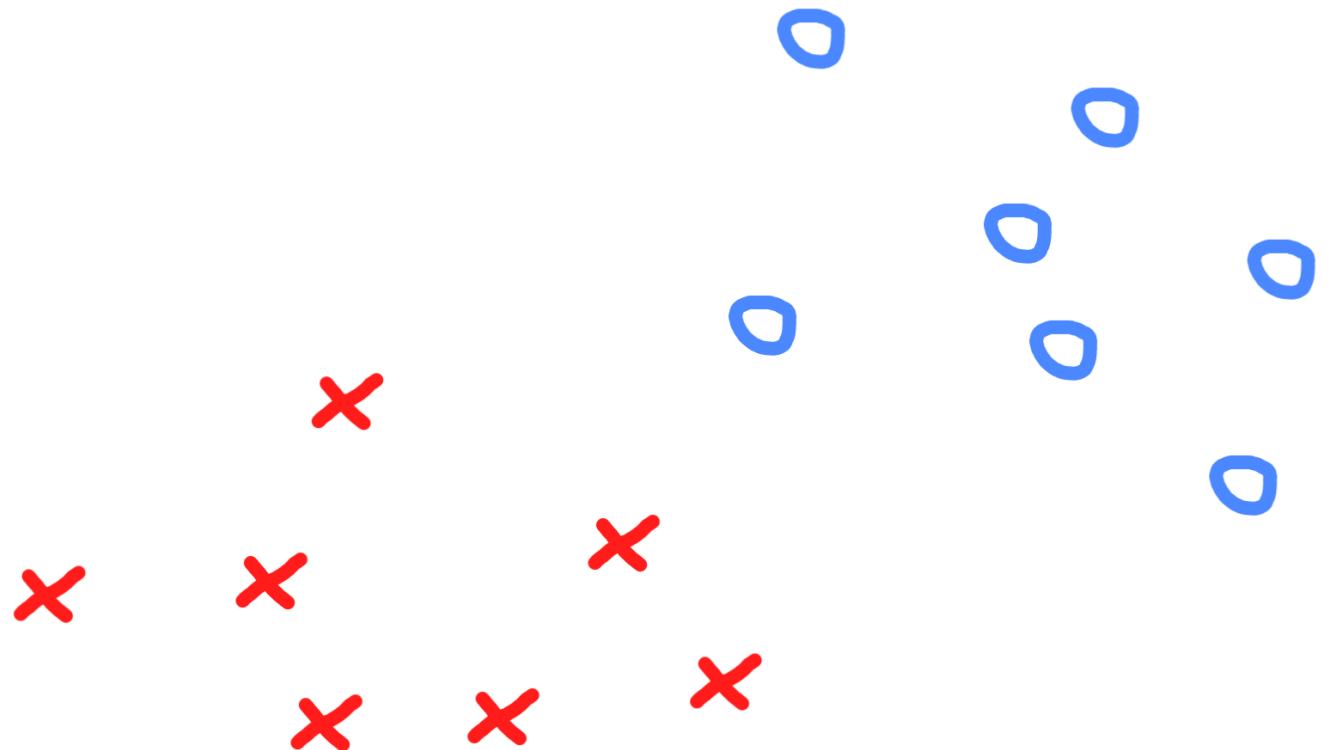


(b) Explanation

*“Thing B is similar to thing A,
and thing A caused Y,
so I predict that B will cause Y!”*

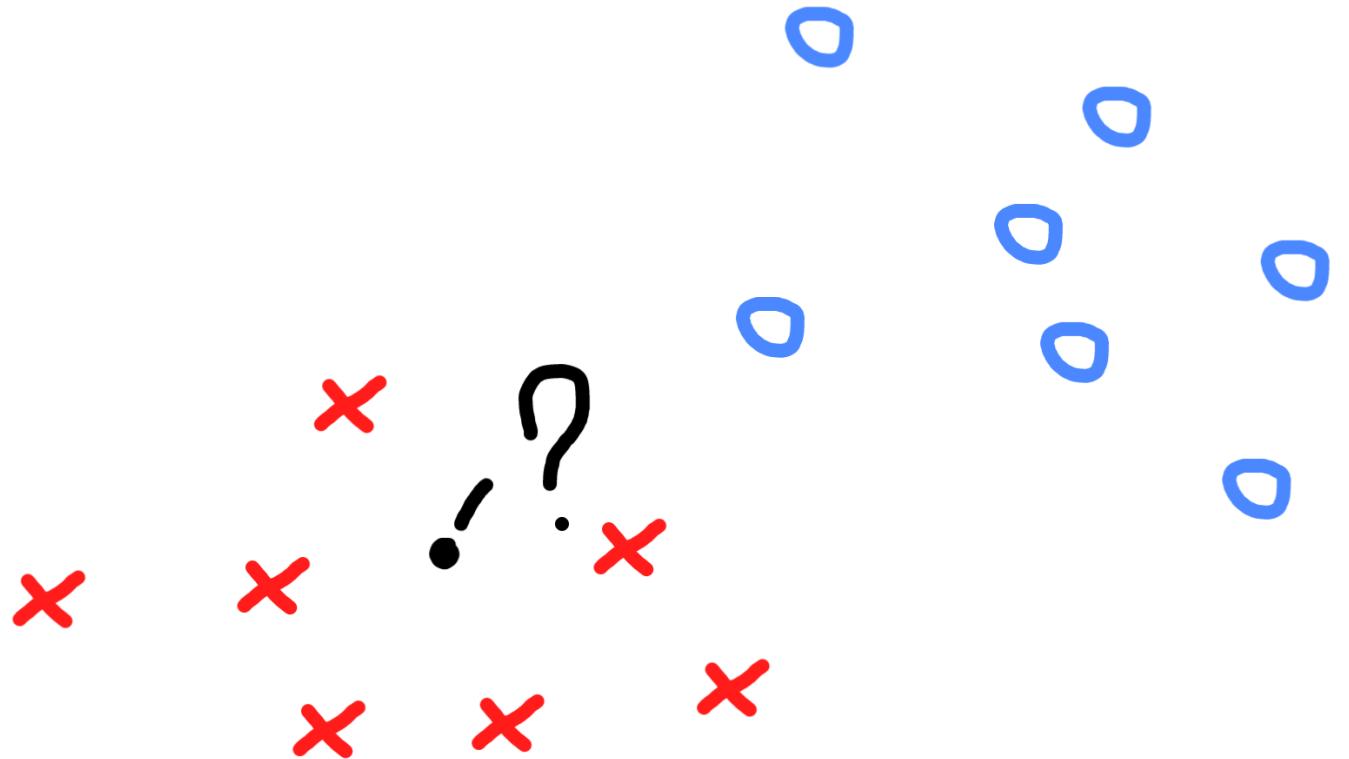
Example-based explanations

1. Some dataset with two classes



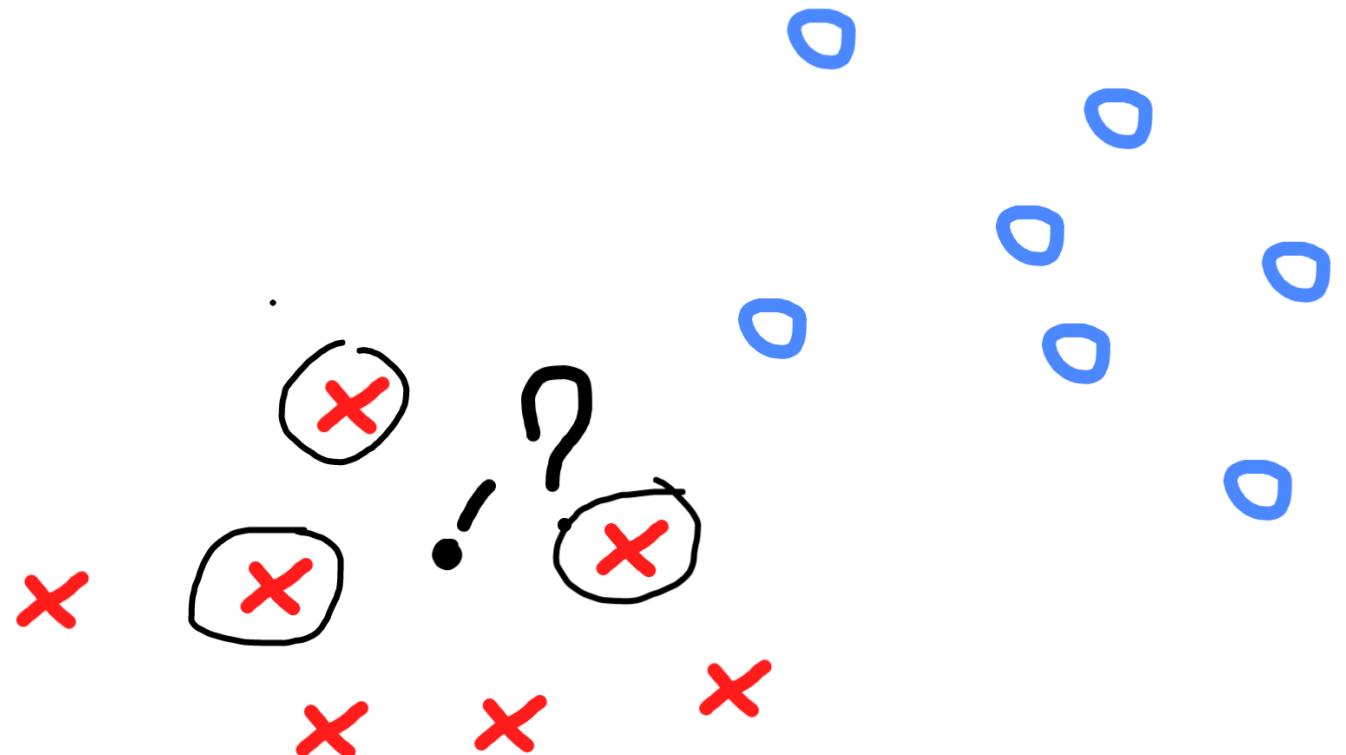
Example-based explanations

1. Some dataset with two classes
2. A new data point to classify



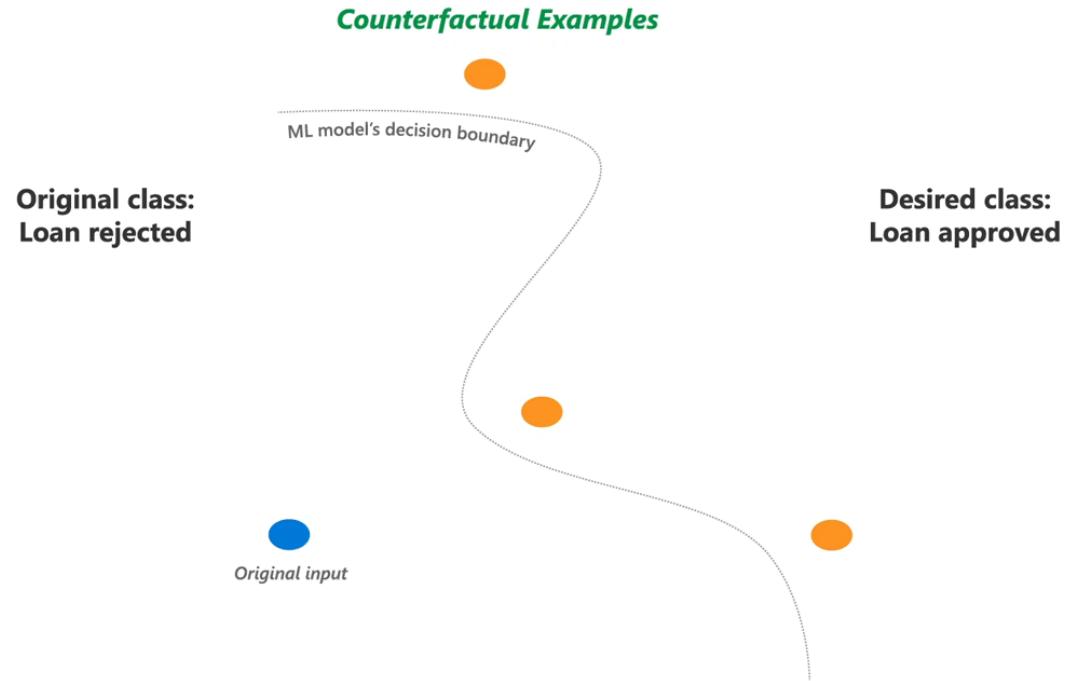
Example-based explanations

1. Some dataset with two classes
2. A new data point to classify
3. An explanation based on nearby data points



Example-based explanations

- Some models
- Counterfactual explanations
- Adversarial examples
- Prototypes & Criticisms
- Influential instances



Jupyter Notebook 1

01_Introduction.ipynb

The inherent interpretability of ML models

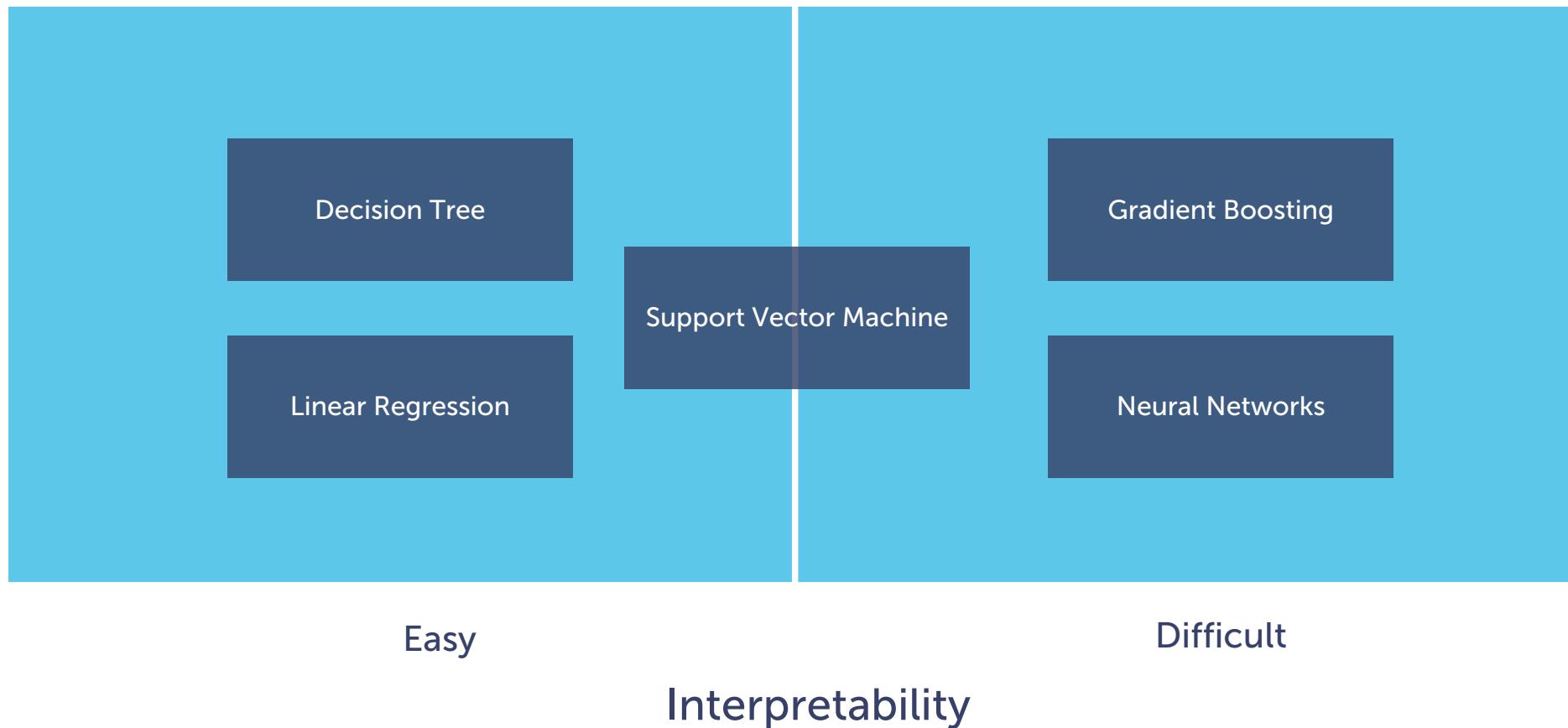
Penguin 218

- Species: Chinstrap
- Flipper length: 210mm
- Body mass: 4100g

MISCLASSIFIED AS GENTOO

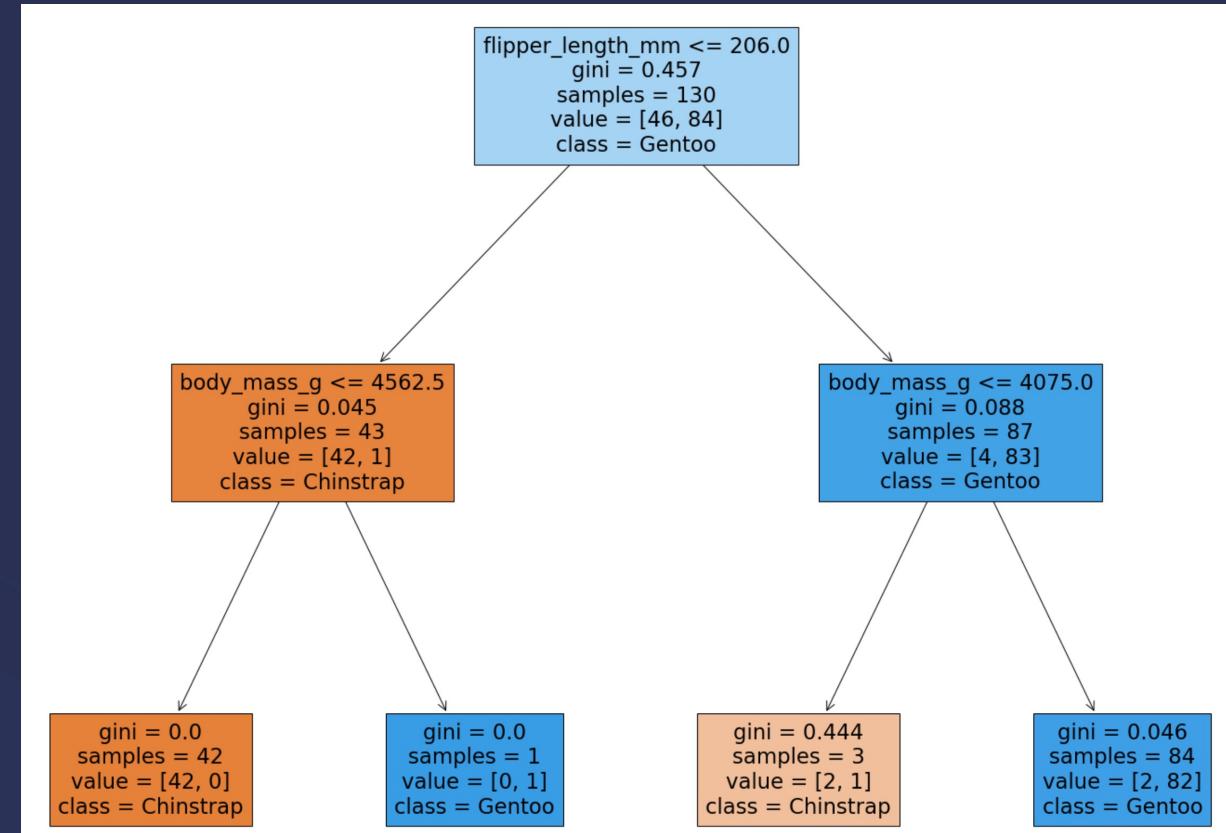


Inherent model interpretability

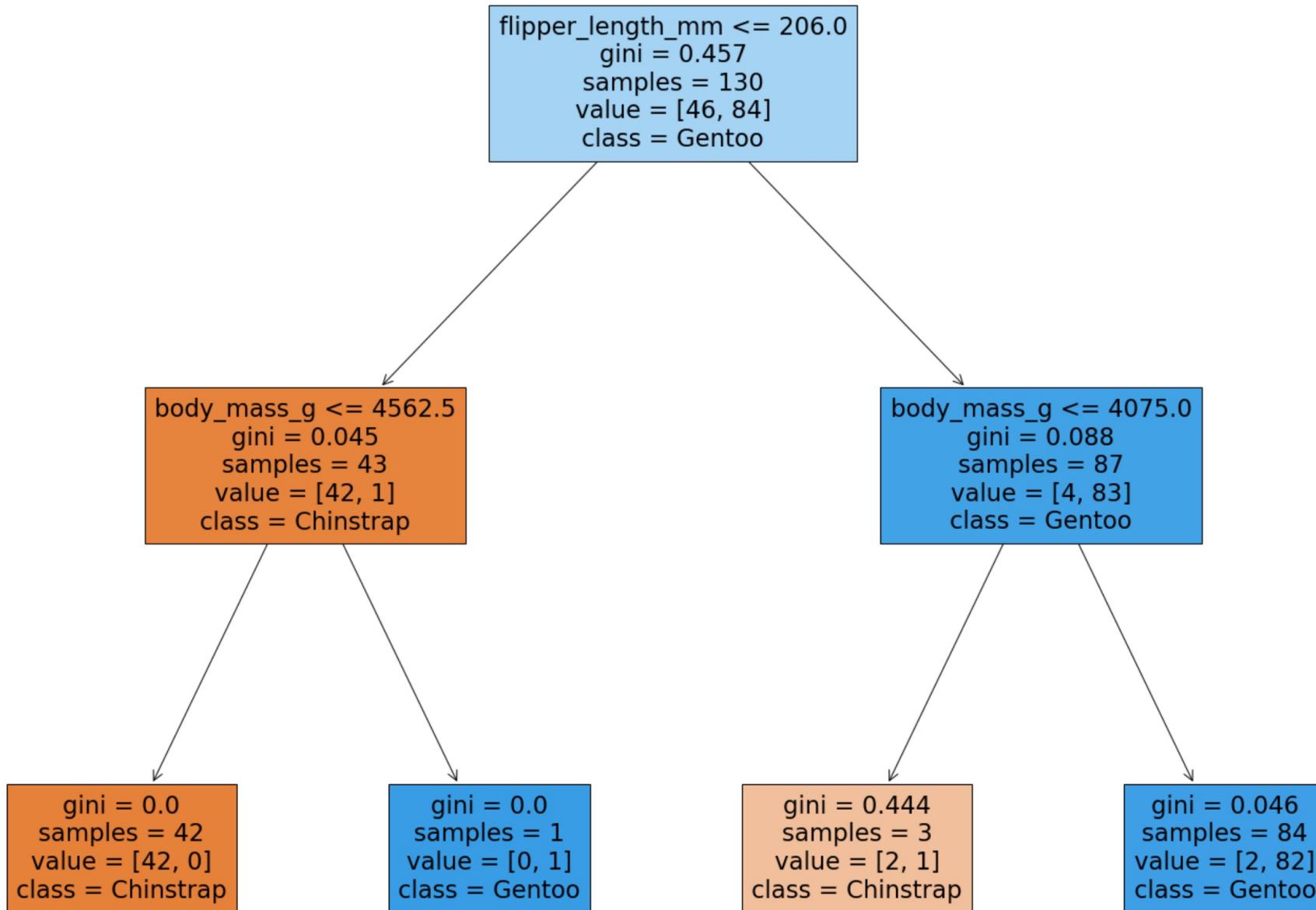


Decision Tree Classifier

The inherent
interpretability of ML
models



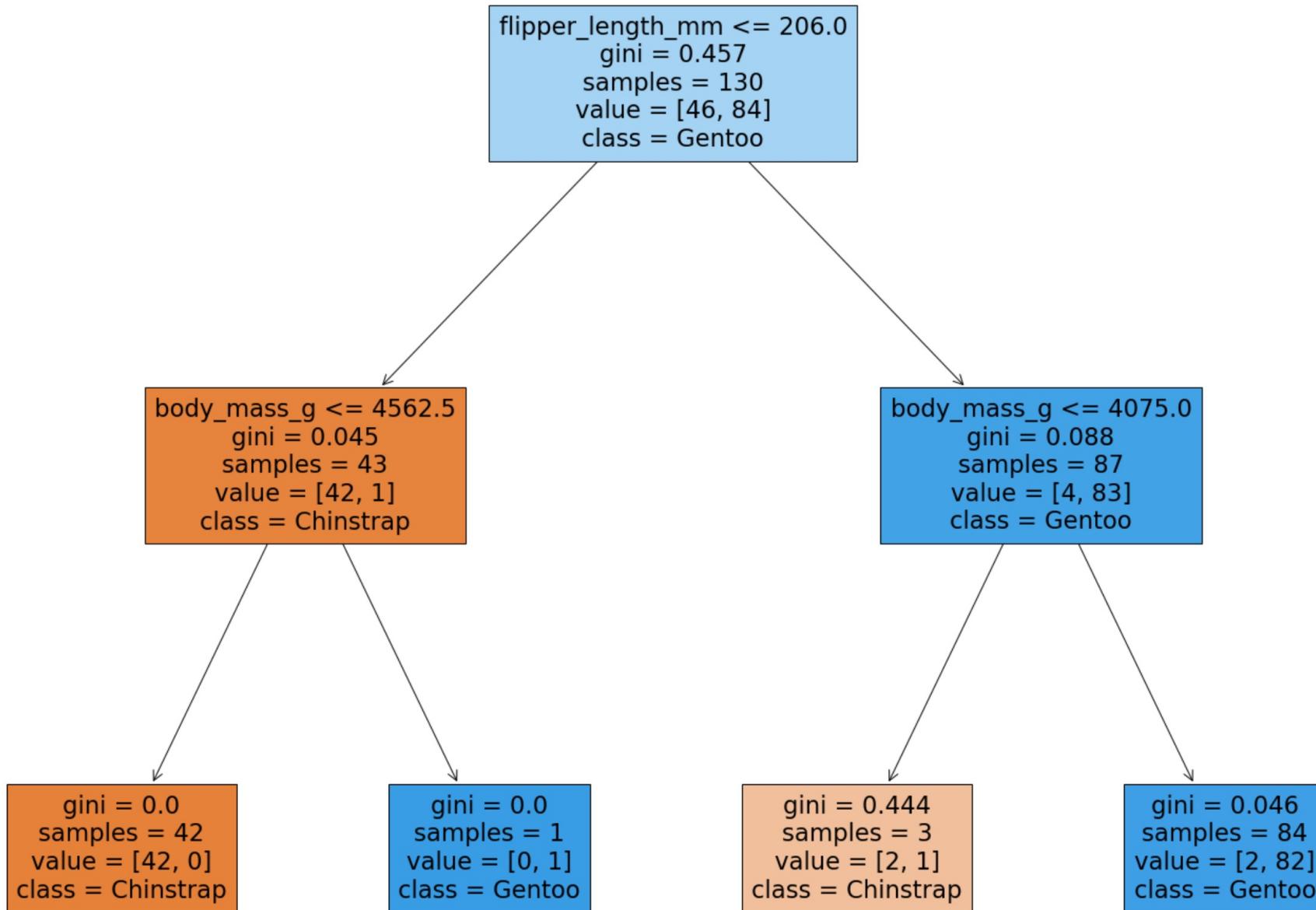
Decision Tree Classifier



Decision Tree Classifier



Decision Tree Classifier

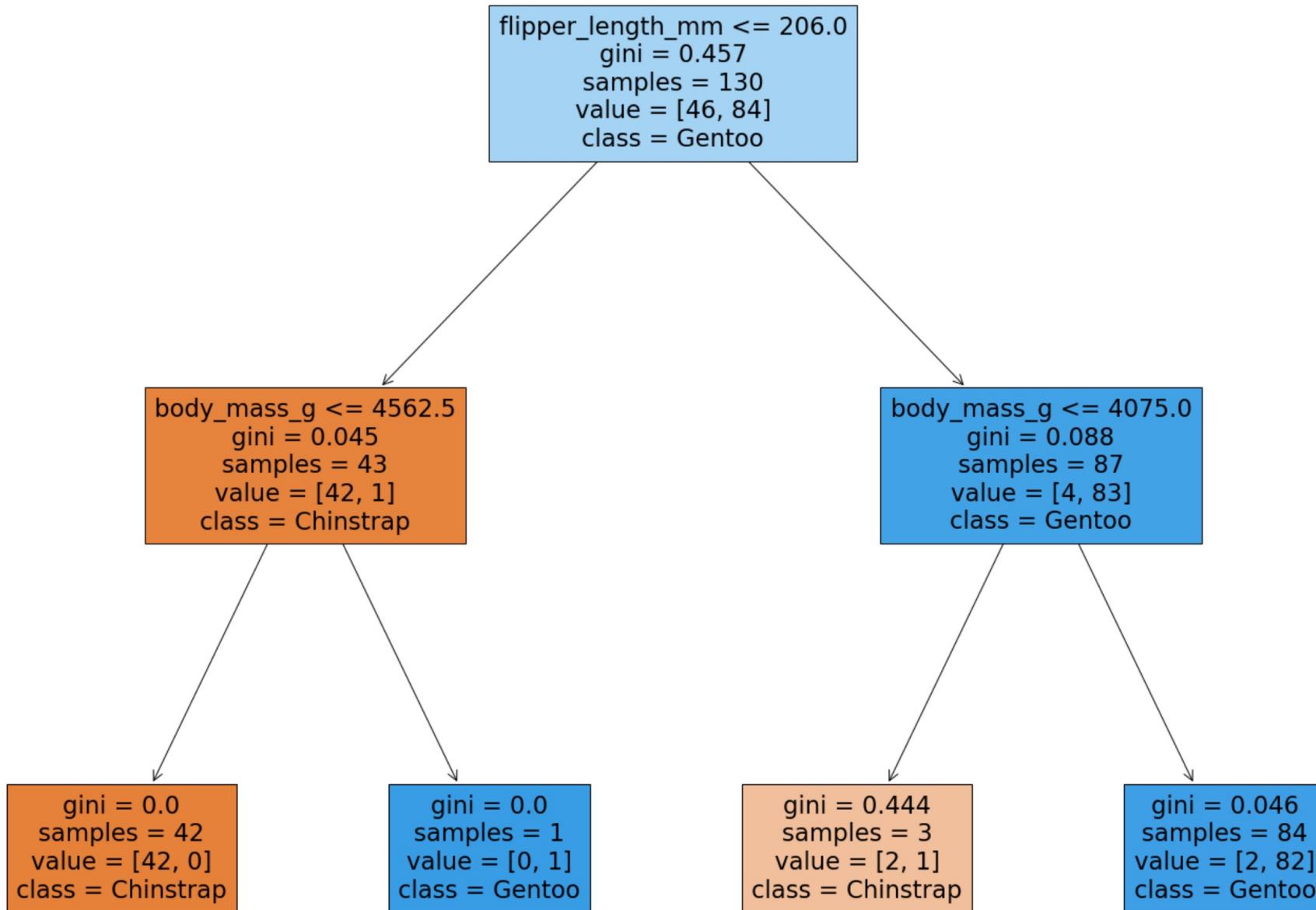


Decision Tree Classifier

```
gini = 0.457
samples = 130
value = [46, 84]
class = Gentoo
```

$$\text{gini} = 1 - \frac{46^2}{130} - \frac{84^2}{130} = 0.457\dots$$

Decision Tree Classifier

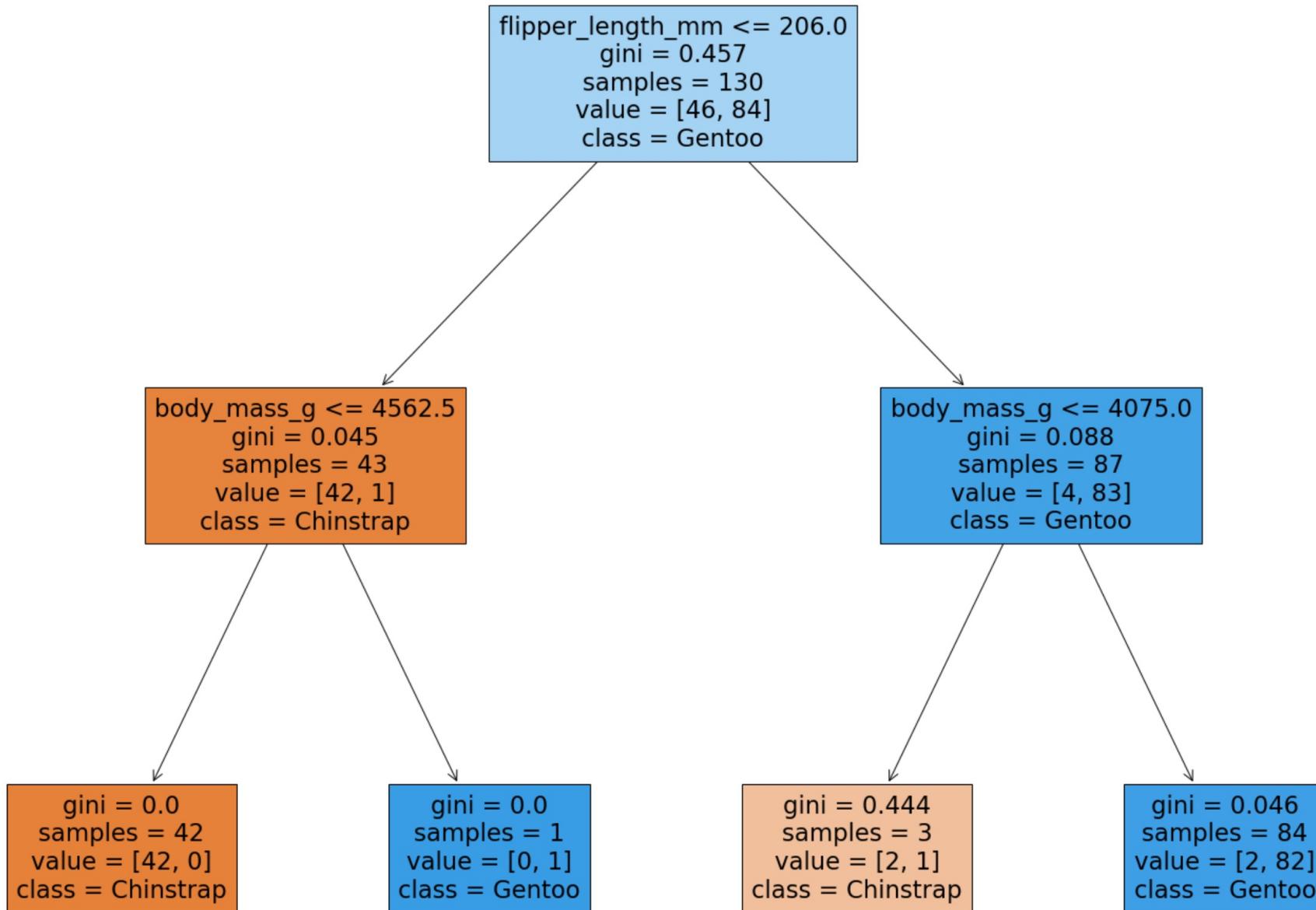


Decision Tree Classifier

```
gini = 0.0
samples = 42
value = [42, 0]
class = Chinstrap
```

$$\text{gini} = 1 - \frac{42^2}{42} - \frac{0^2}{42} = 1 - 1 - 0 = 1$$

Decision Tree Classifier

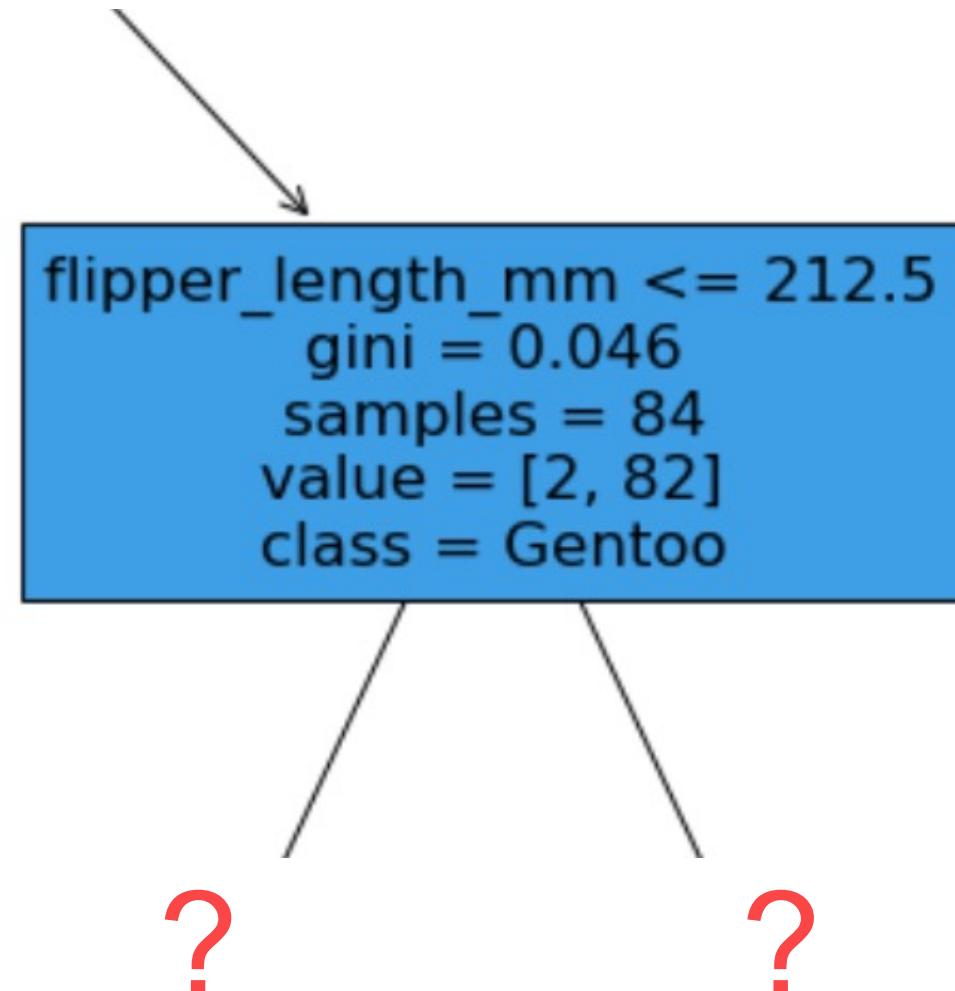


Decision Tree Classifier

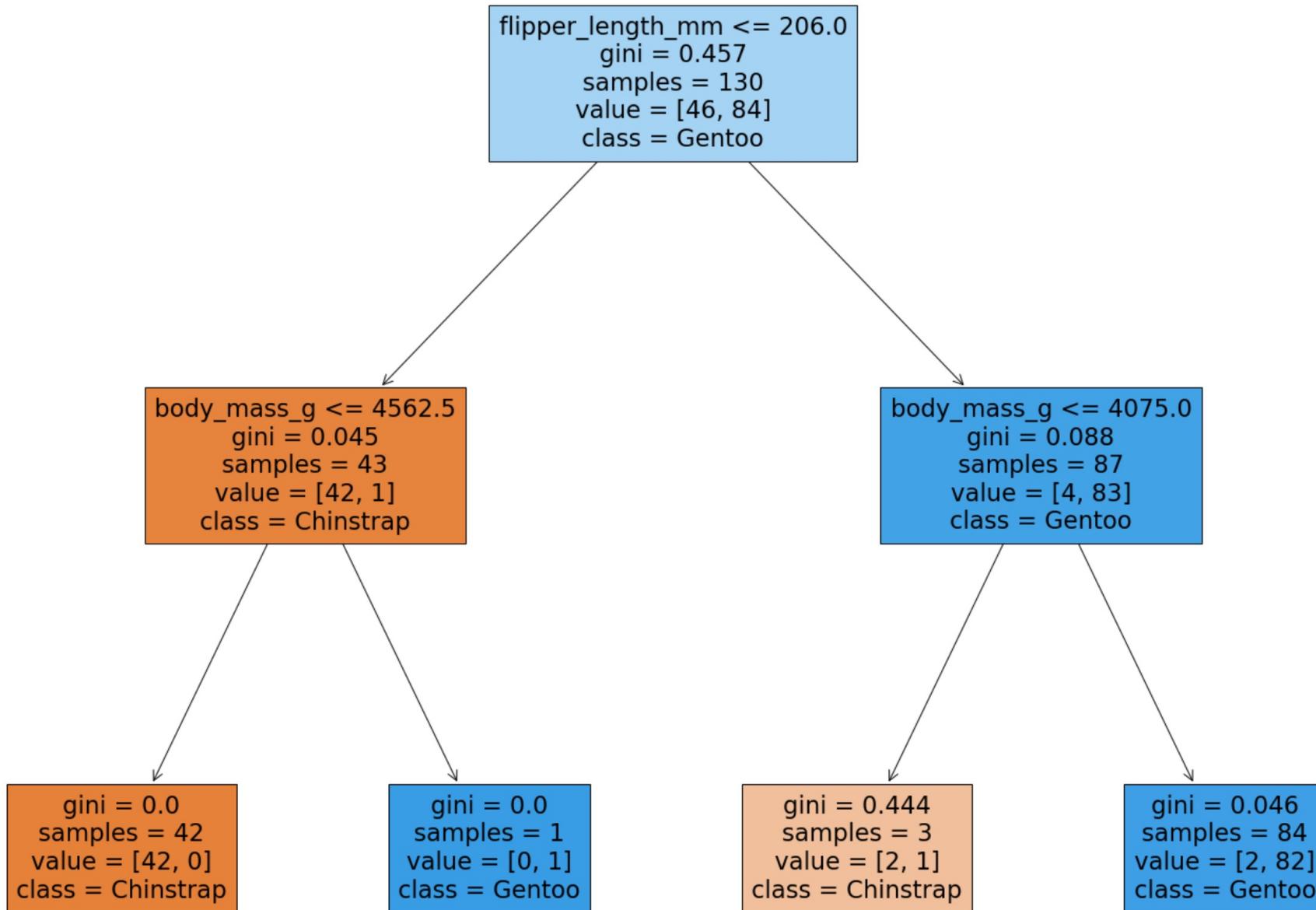
```
gini = 0.046
samples = 84
value = [2, 82]
class = Gentoo
```

$$\text{gini} = 1 - \frac{2^2}{84} - \frac{82^2}{84} = 0.046$$

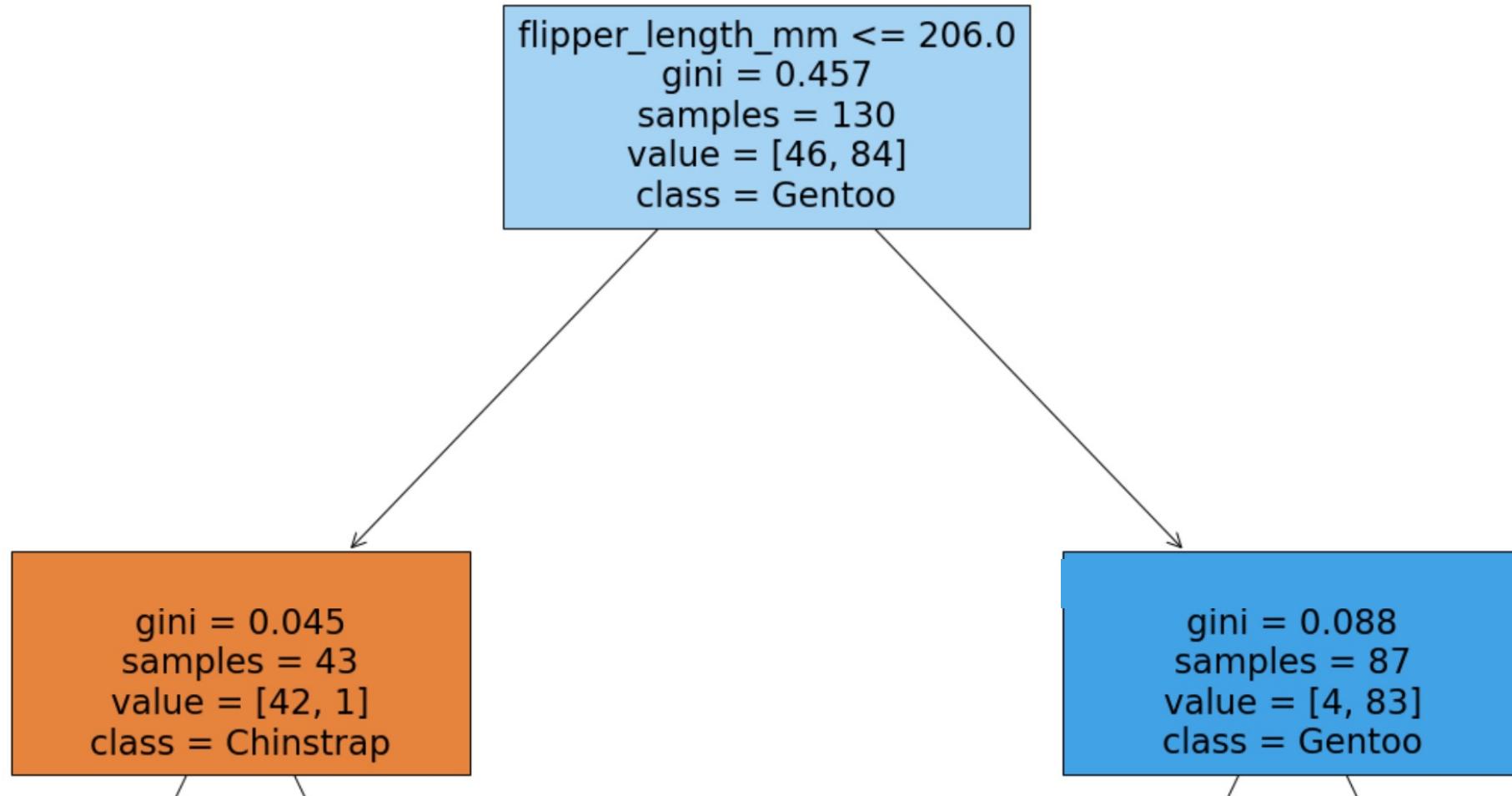
Decision Tree Classifier



Decision Tree Classifier



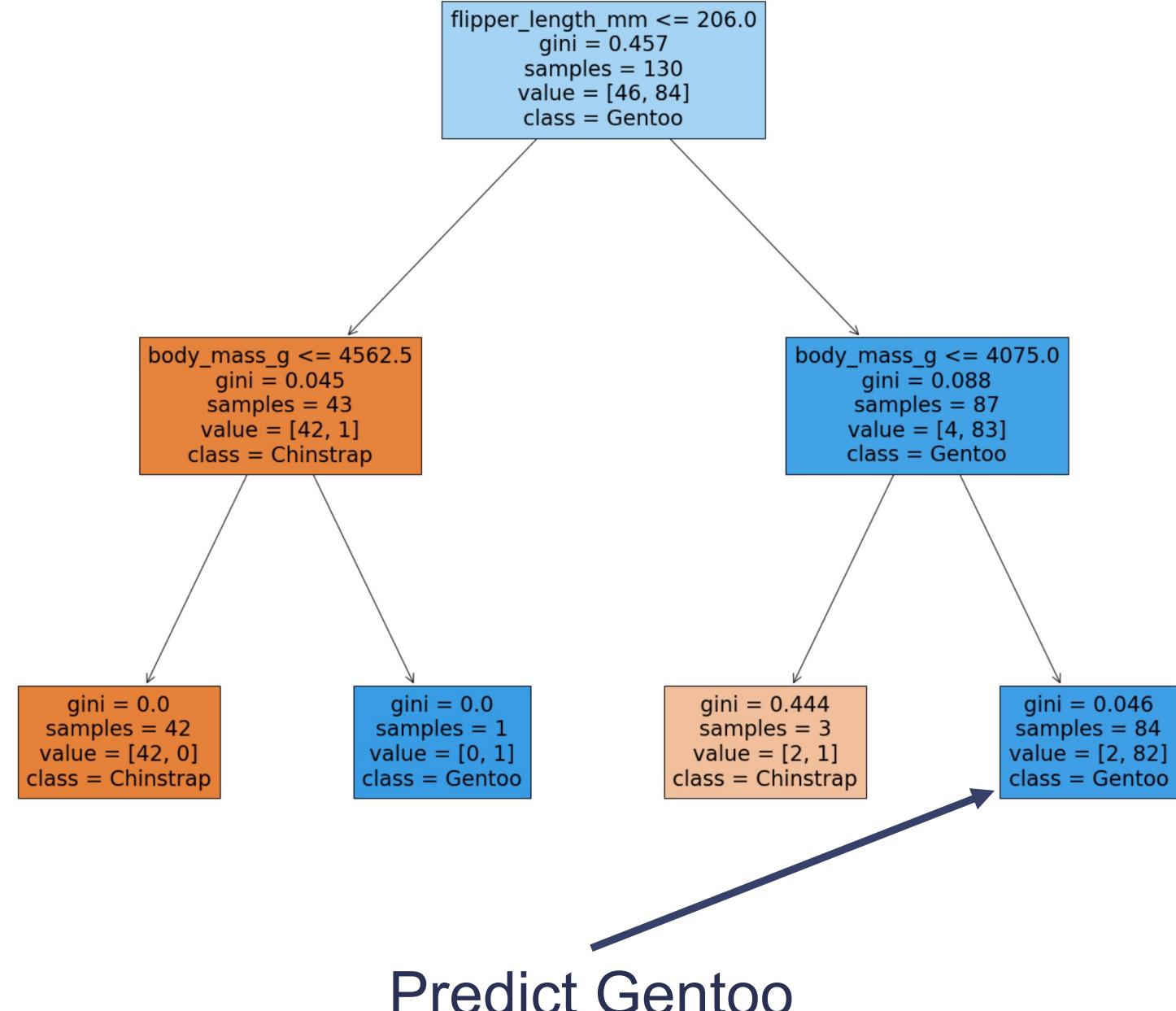
Decision Tree Classifier



$$0.045 \times \frac{43}{130} + 0.088 \times \frac{87}{130} = 0.074$$

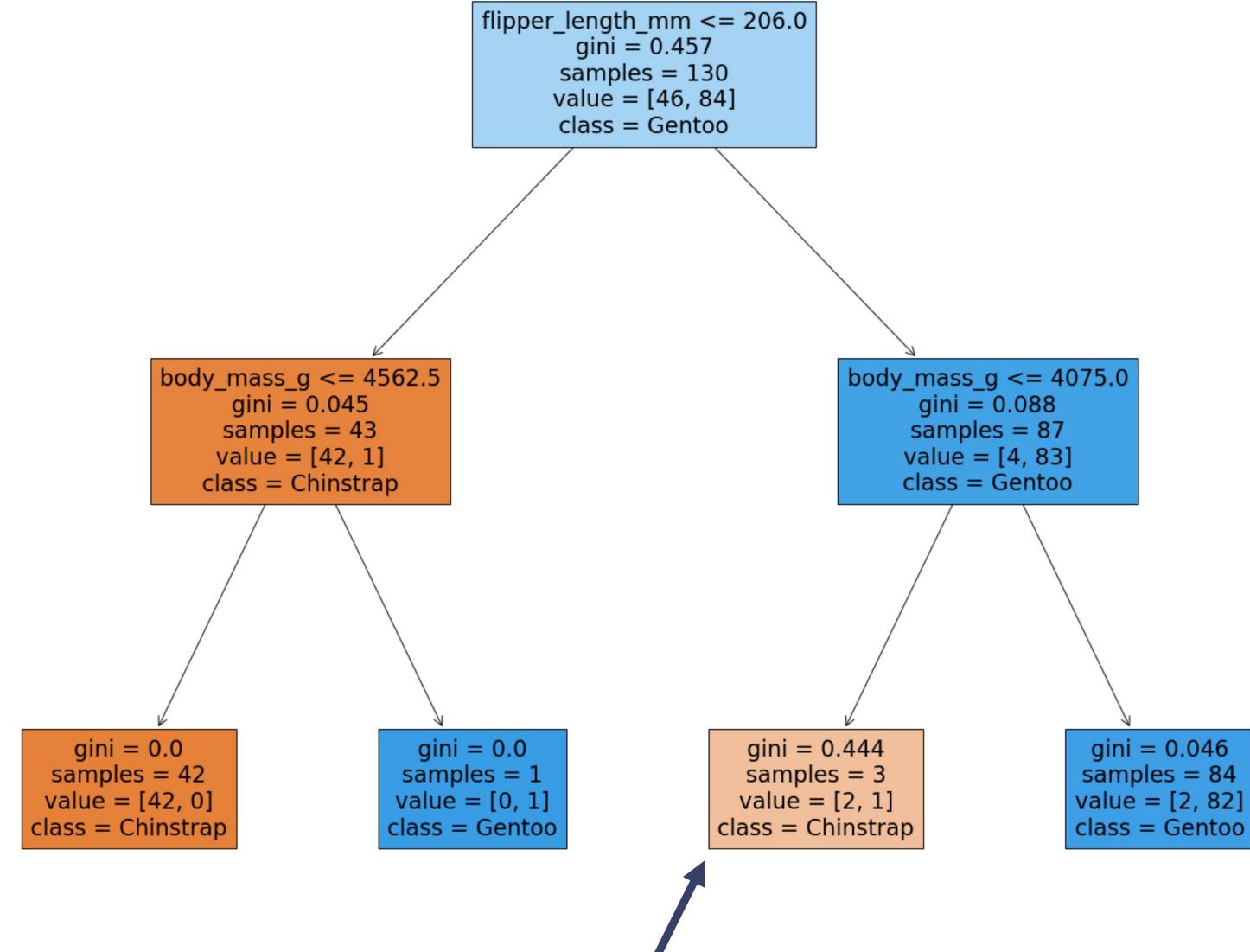
Penguin 218

- Flipper length: 210mm
- Body mass: 4100g



Penguin 218

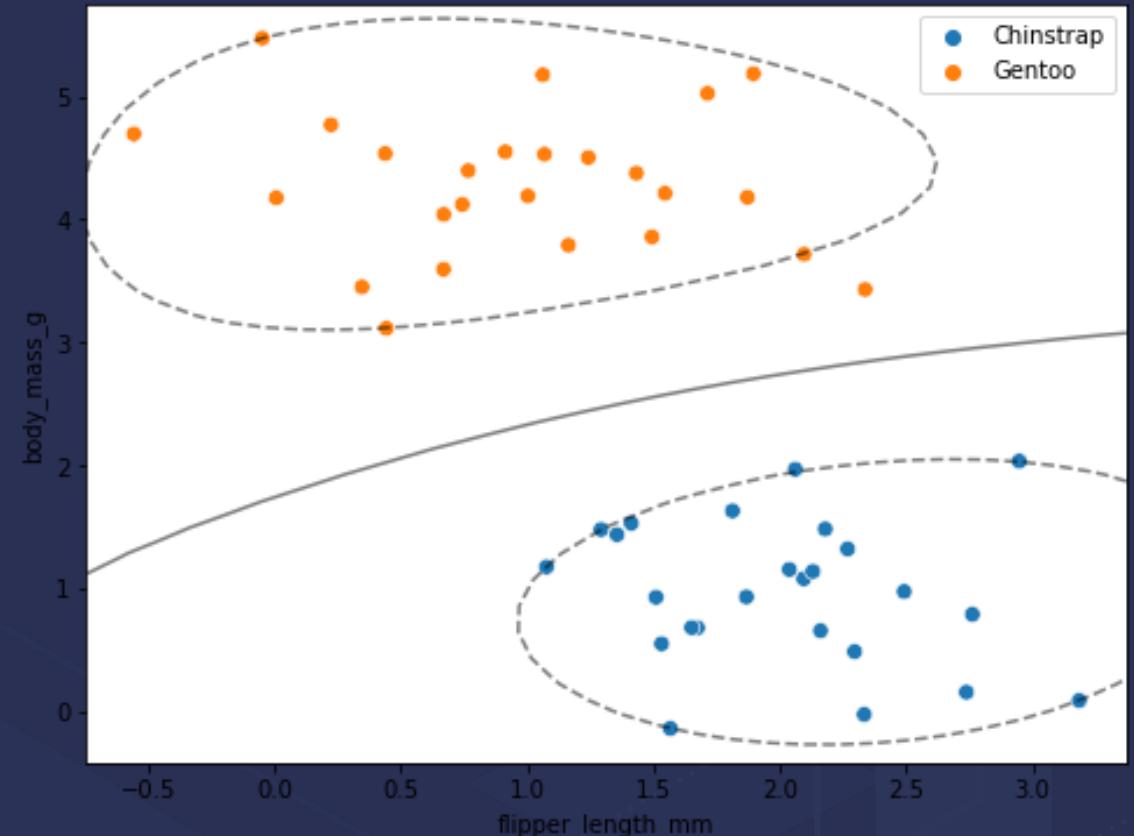
- Flipper length: 210mm
- Body mass: 4070g



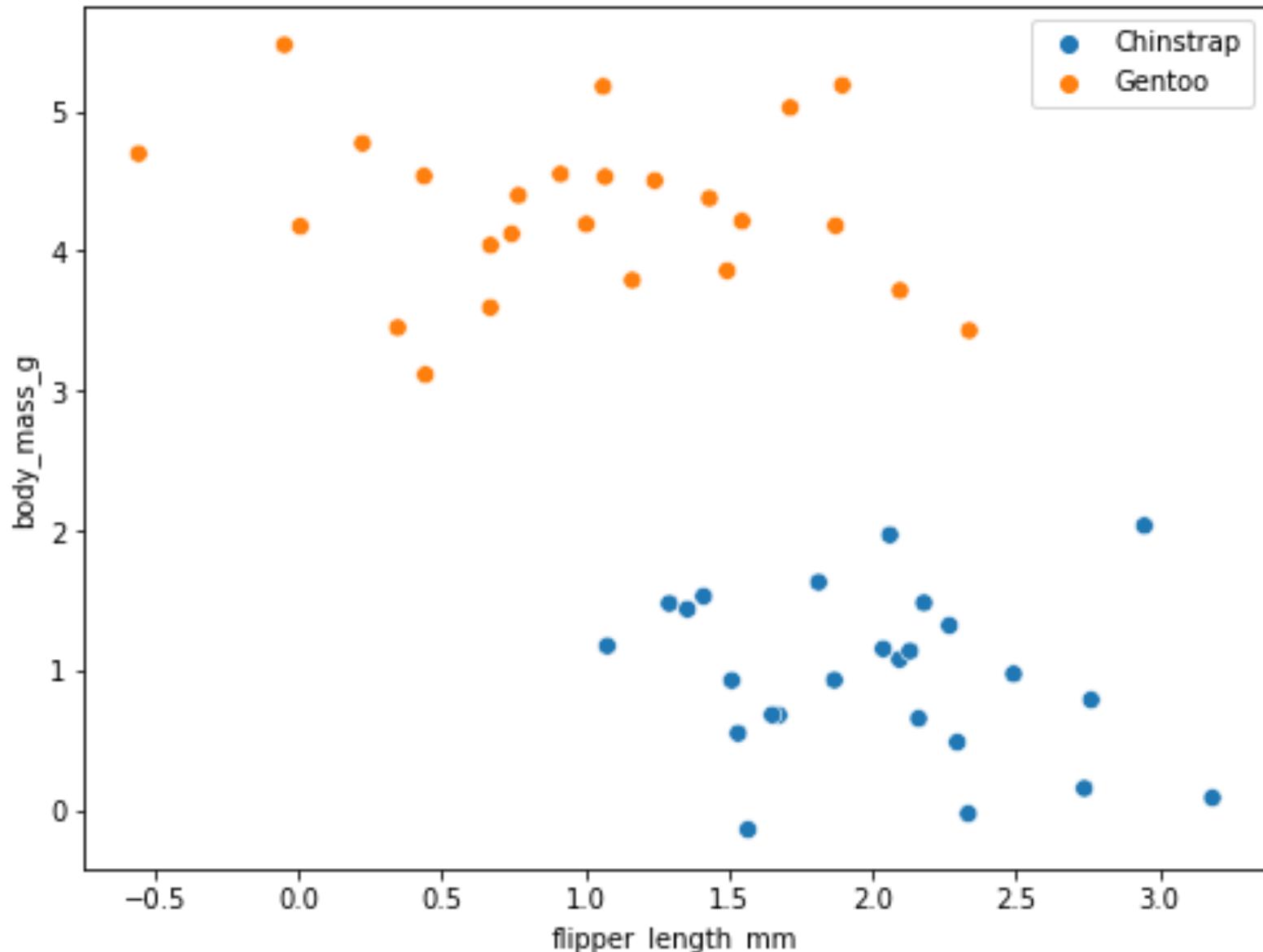
Predict Chinstrap

Support Vector Machine Classifier

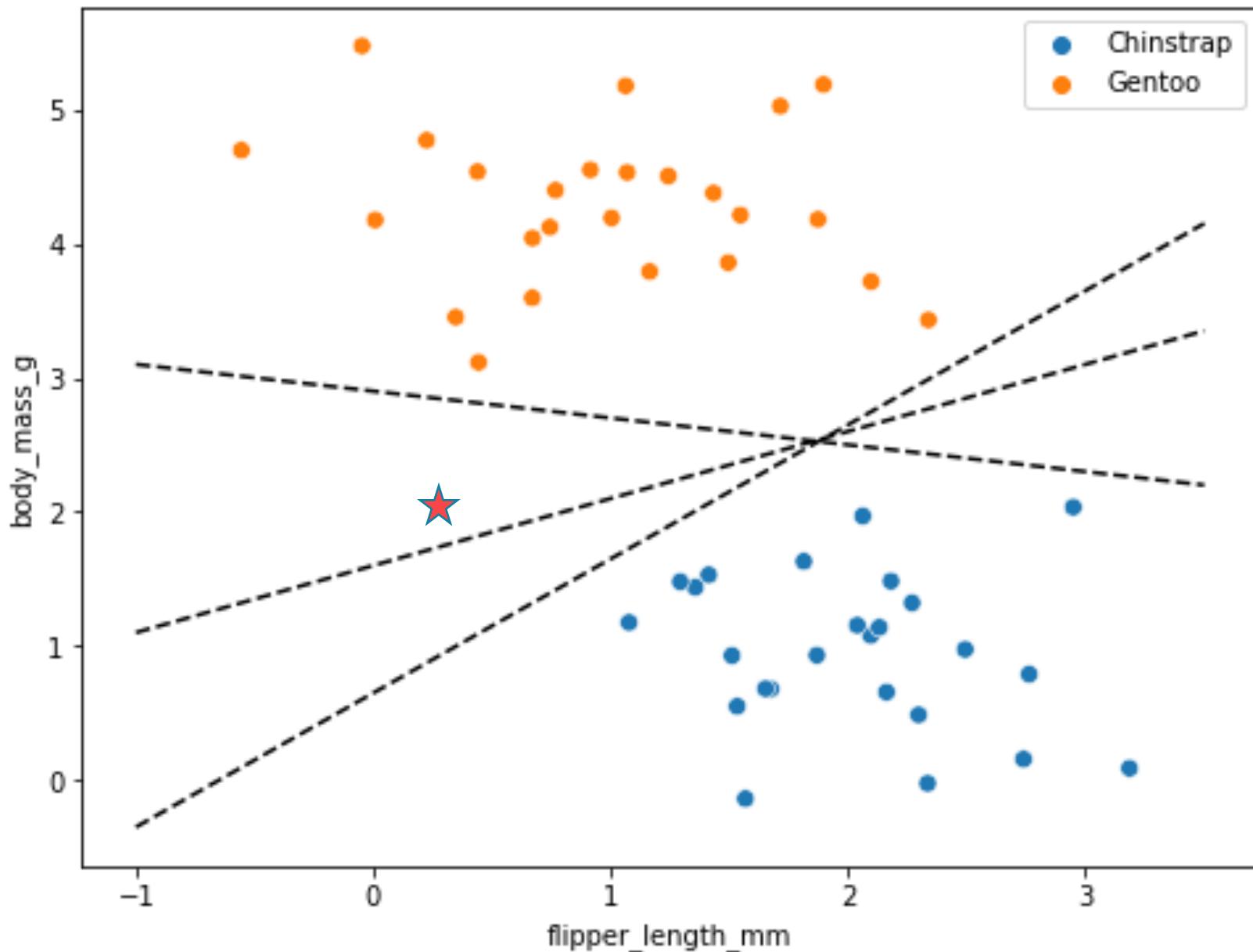
The inherent
interpretability of ML
models



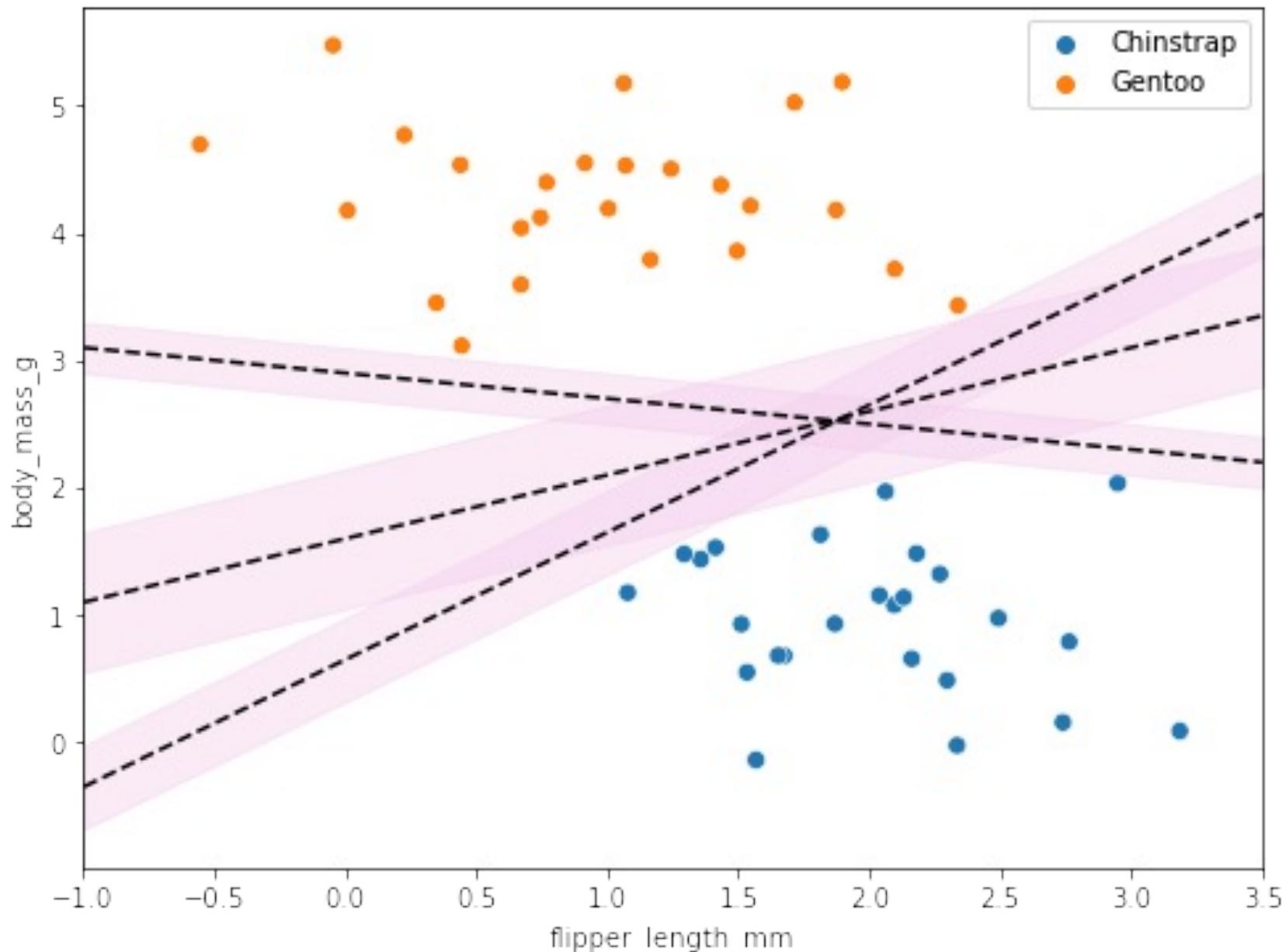
Support Vector Machine Classifier



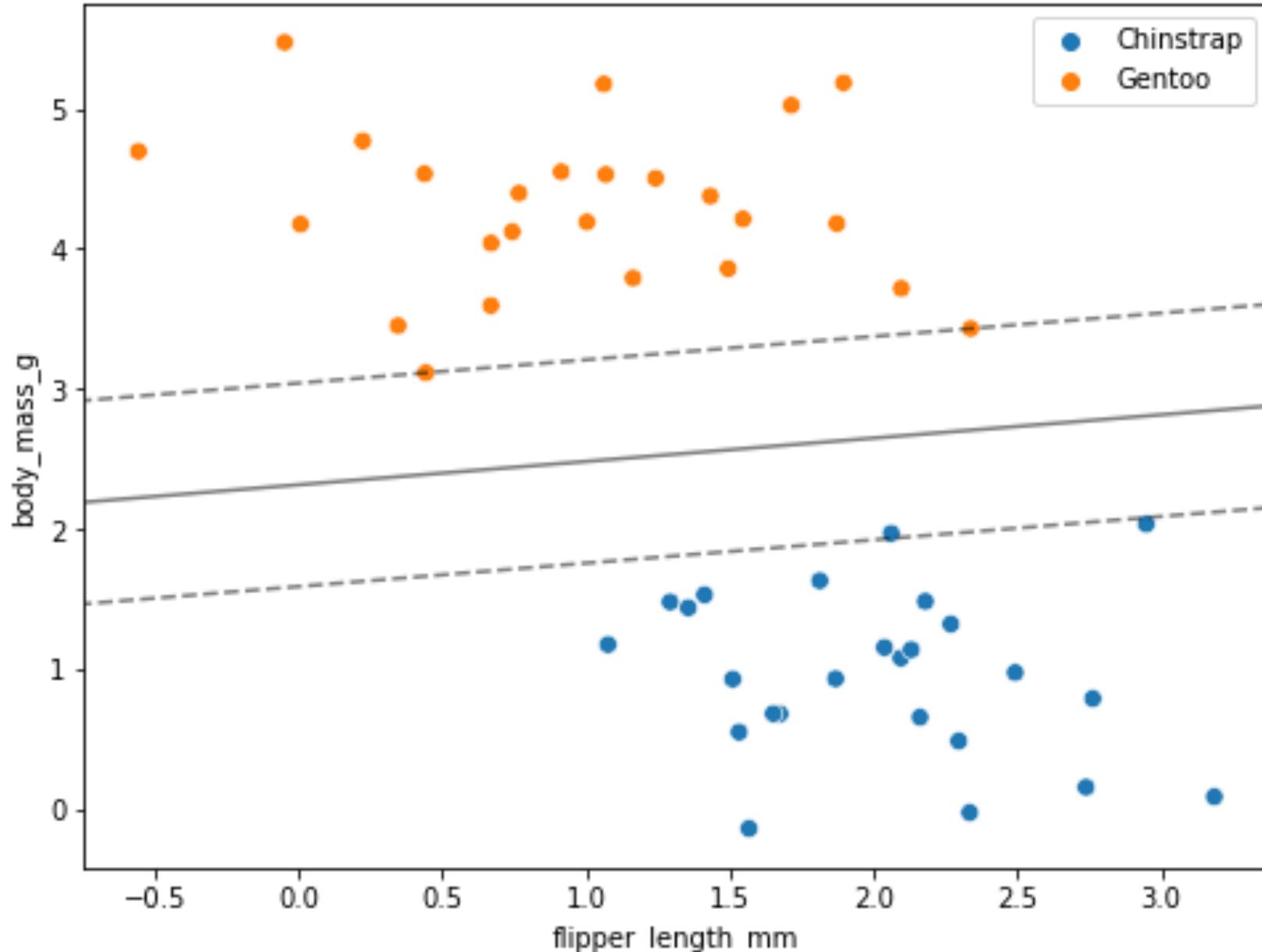
Support Vector Machine Classifier



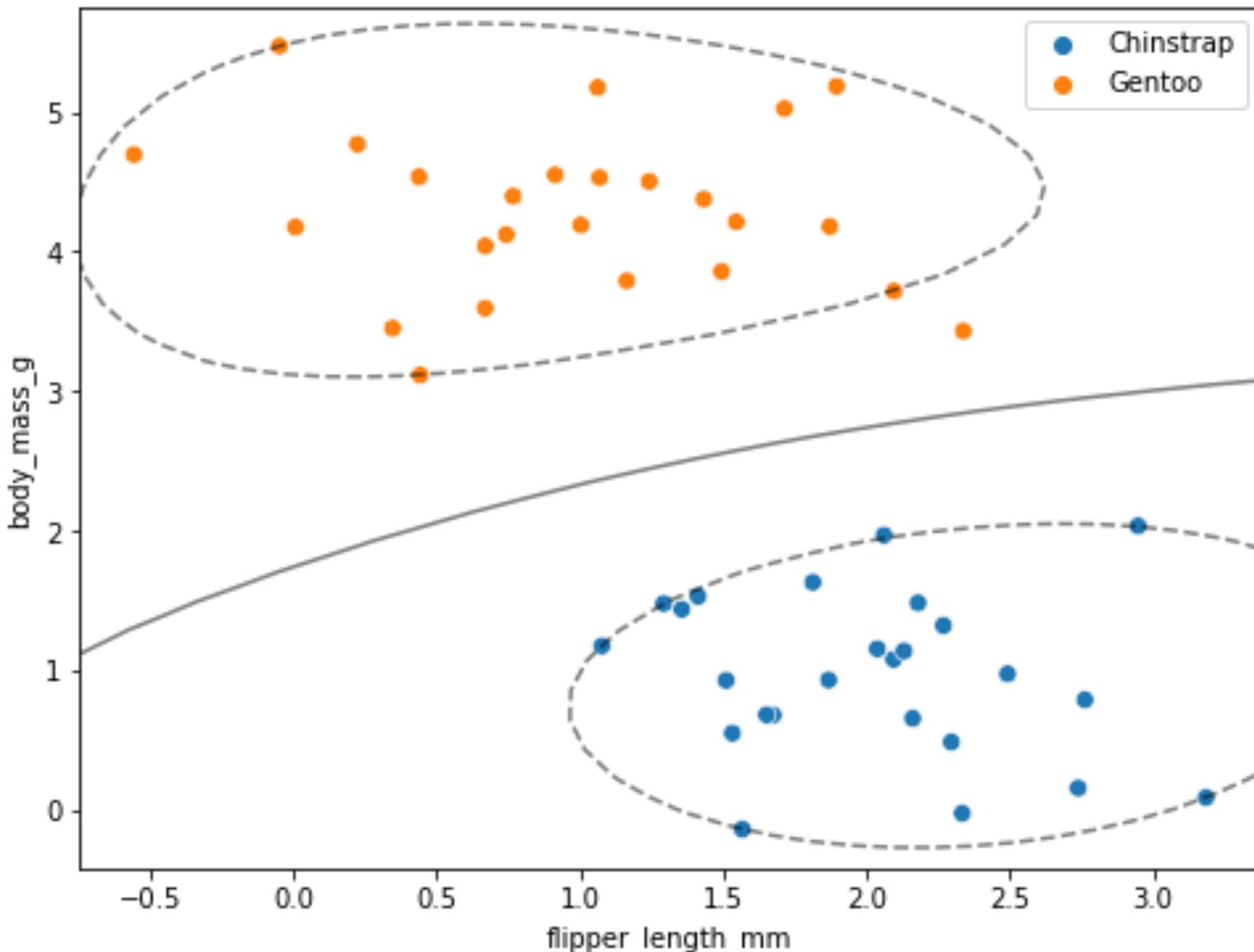
Support Vector Machine Classifier



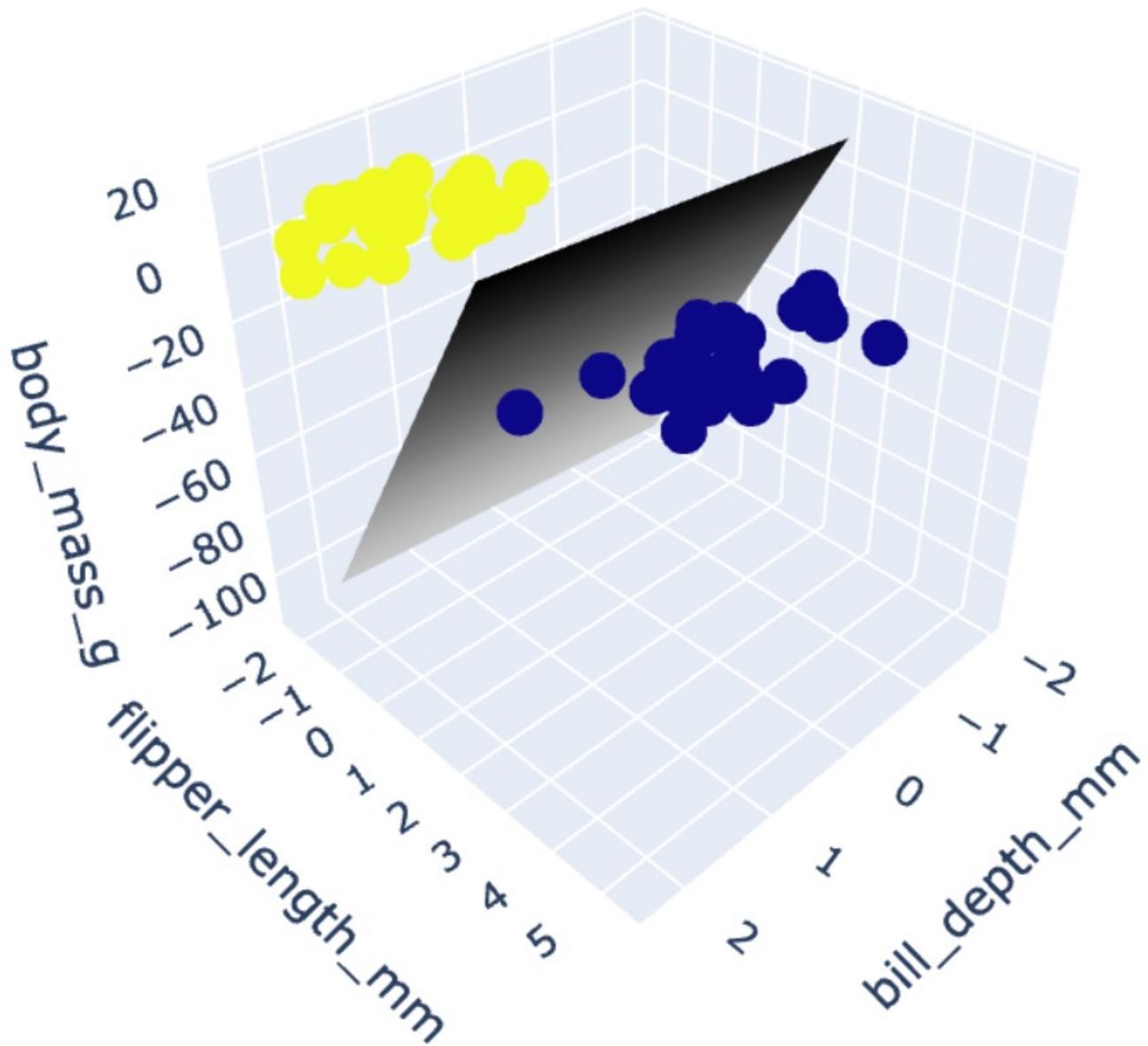
Support Vector Machine Classifier



Support Vector Machine Classifier

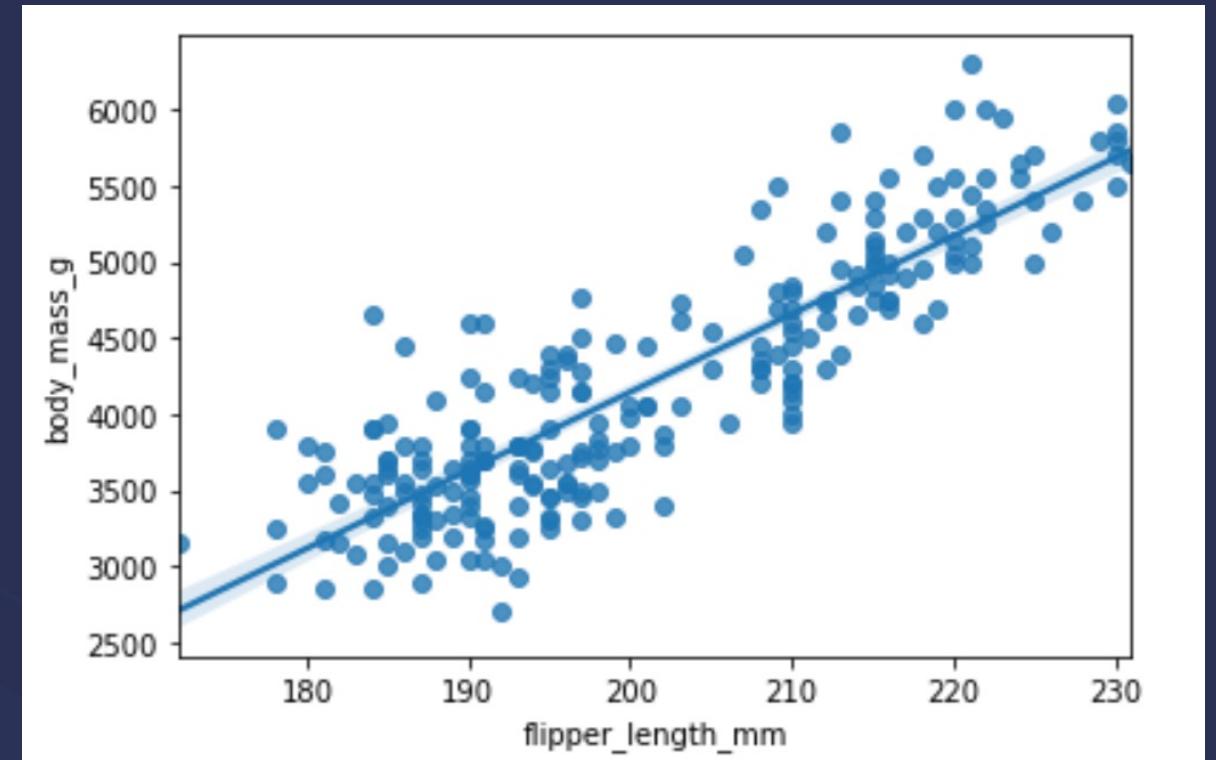


Support Vector Machine Classifier



Linear Regression

The inherent
interpretability of ML
models



Classification

bill_depth_mm	flipper_length_mm	species
18.7	181.0	Adelie
17.4	186.0	Adelie
18.0	195.0	Adelie
19.3	193.0	Adelie
20.6	190.0	Adelie



Regression

bill_depth_mm	flipper_length_mm	body_mass_g
18.7	181.0	3750.0
17.4	186.0	3800.0
18.0	195.0	3250.0
19.3	193.0	3450.0
20.6	190.0	3650.0



Linear Regression

$$m * \text{bill depth} + n * \text{flipper length} + c = \text{body mass}$$

bill_depth_mm	flipper_length_mm	body_mass_g
18.7	181.0	3750.0
17.4	186.0	3800.0
18.0	195.0	3250.0
19.3	193.0	3450.0
20.6	190.0	3650.0



Linear Regression

$m * \text{bill depth} + n * \text{flipper length} + c = \text{body mass}$

$14 * \text{bill depth} + 53 * \text{flipper length} - 6600 = \text{body mass}$

bill_depth_mm	flipper_length_mm	body_mass_g
18.7	181.0	3750.0

$$14 * 18.7 + 53 * 181.0 - 6600 = 3171$$

Model-agnostic interpretability methods

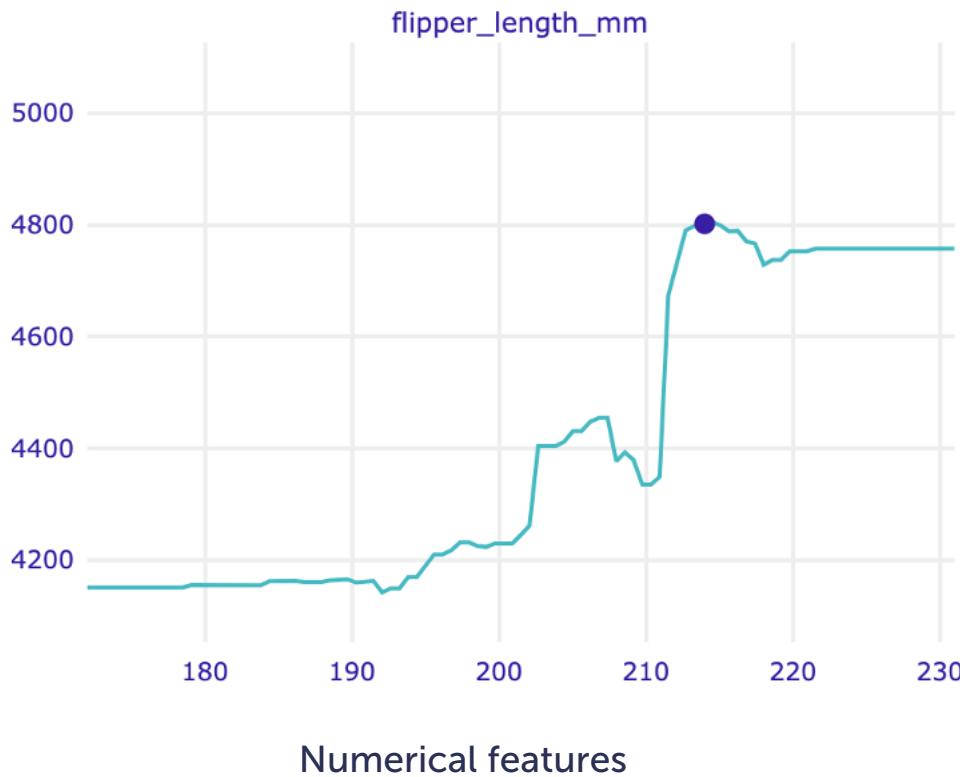
Jupyter Notebook 2

02_Ceteris_Paribus_Plots.ipynb

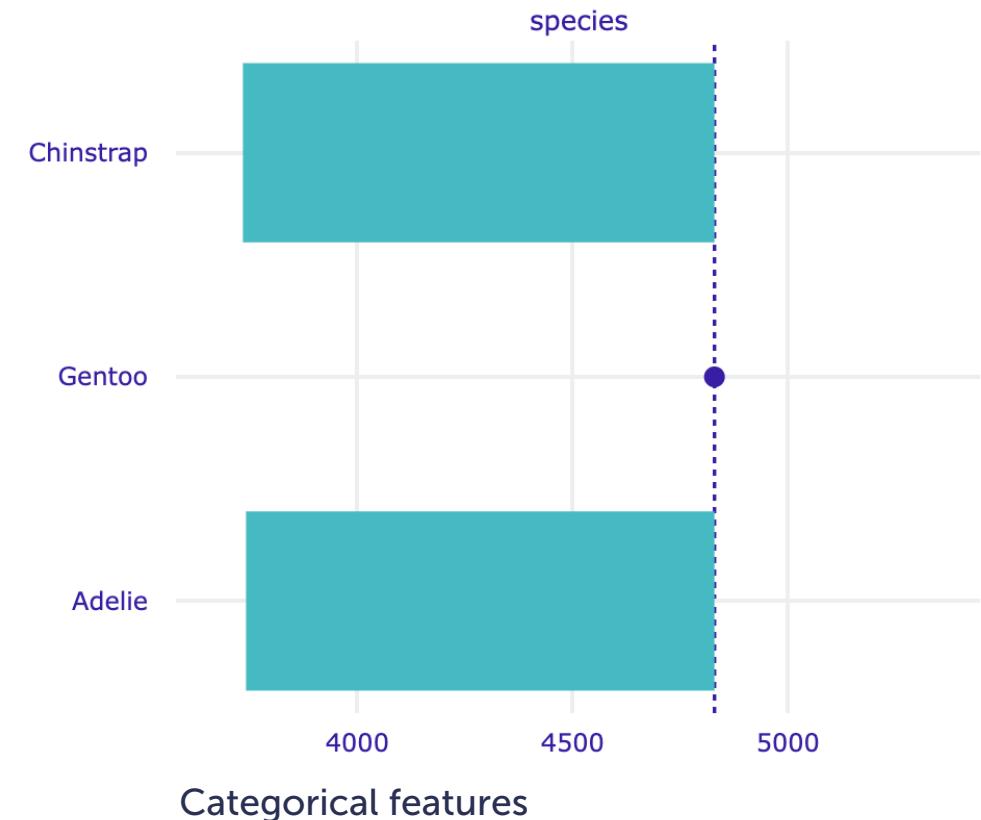
Ceteris Paribus profile

Numerical and categorical features

Ceteris Paribus Profiles



Numerical features

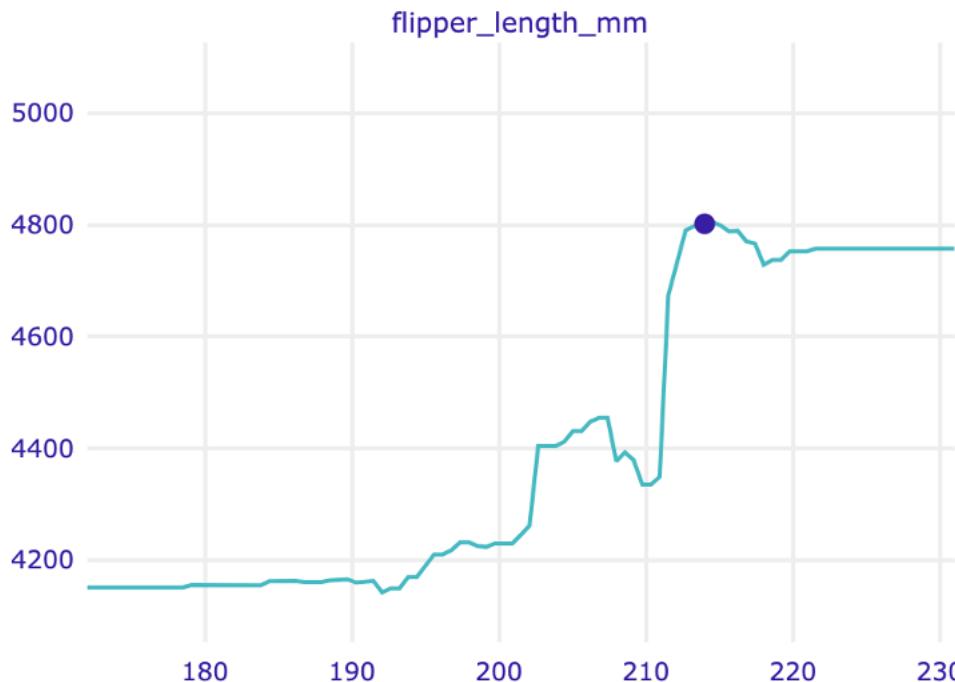


Categorical features

Ceteris Paribus profile

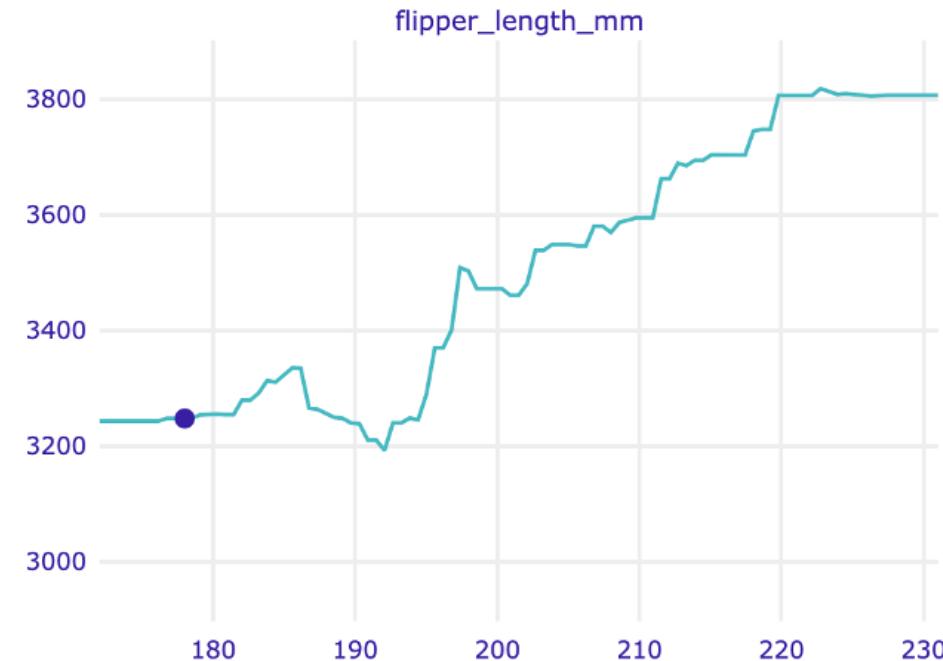
Different plots for different *data points*

Ceteris Paribus Profiles



Row 231: A Gentoo penguin

Ceteris Paribus Profiles

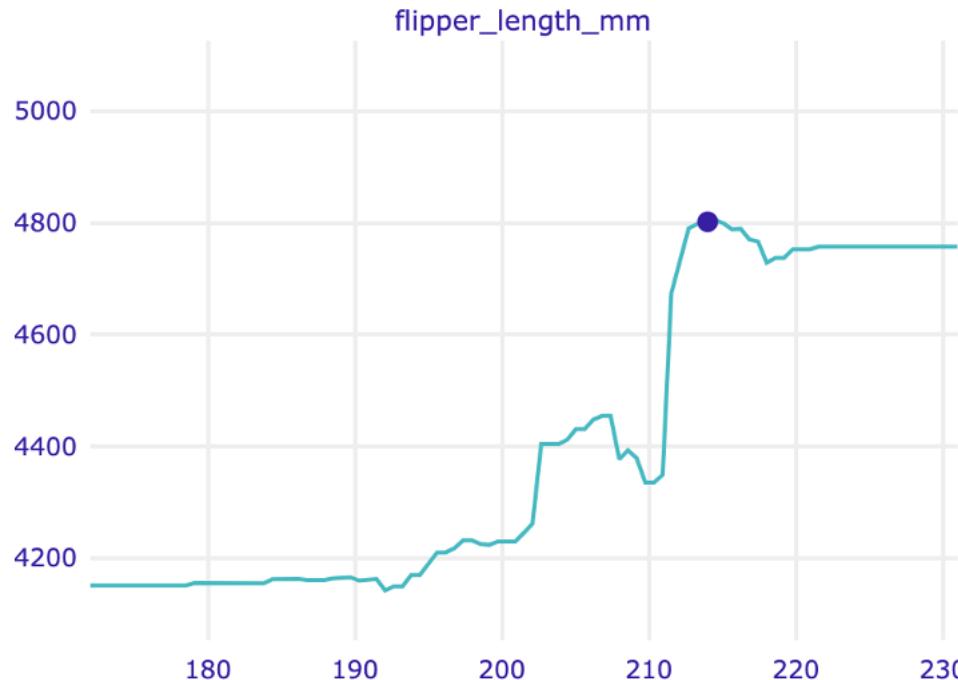


Row 30: An Adélie penguin

Ceteris Paribus profile

Different plots for different *models*

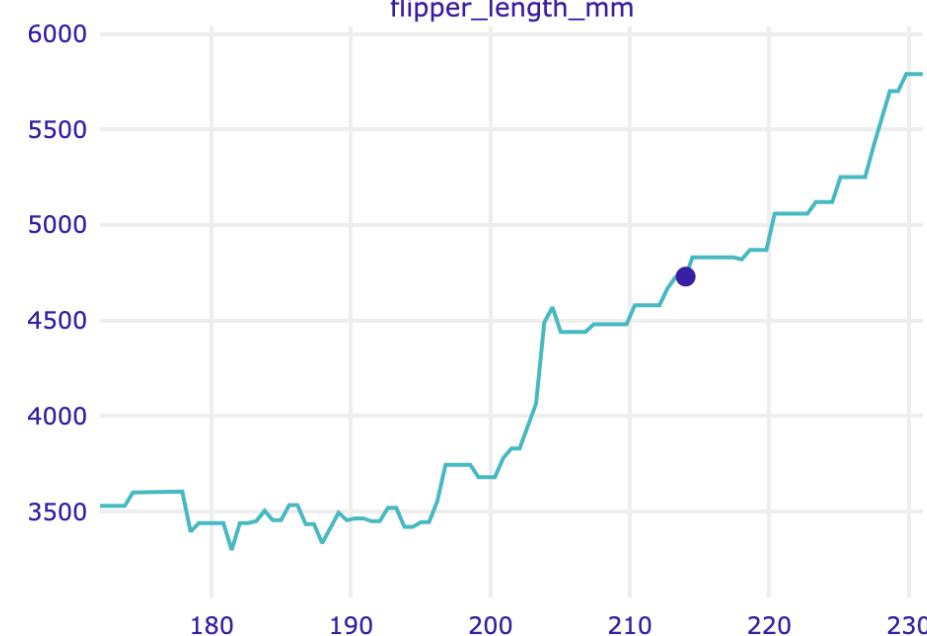
Ceteris Paribus Profiles



Row 231: A Gentoo penguin

Random Forest Regressor

Ceteris Paribus Profiles



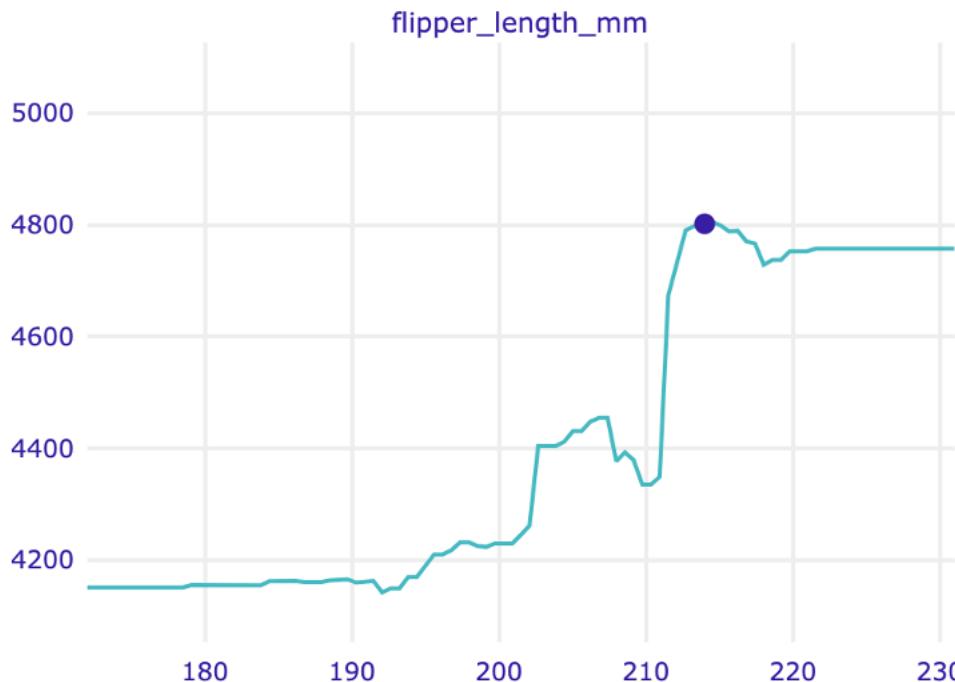
Row 231: A Gentoo penguin

K-Nearest Neighbors Regressor

Ceteris Paribus profile

Different plots for different *models*

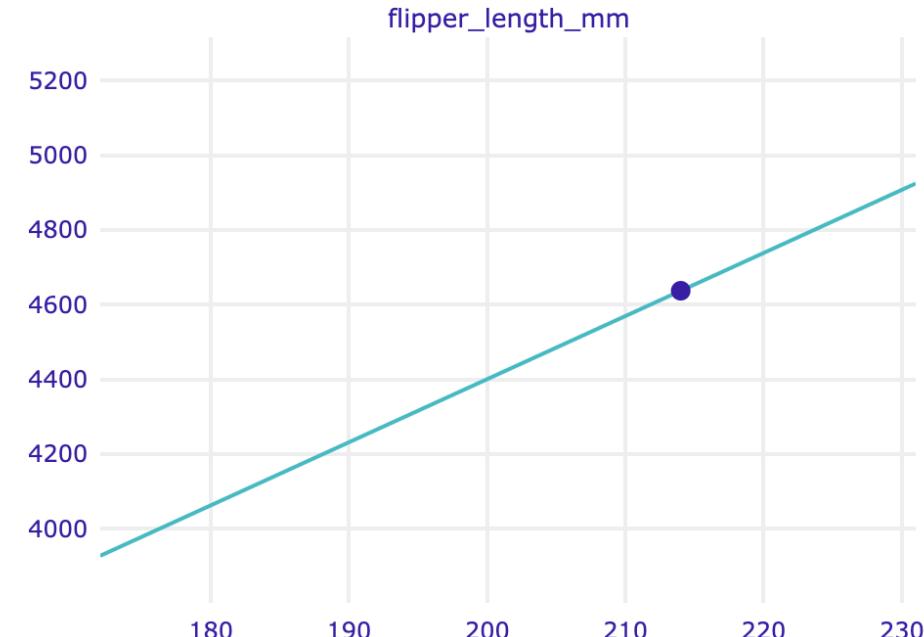
Ceteris Paribus Profiles



Row 231: A Gentoo penguin

Random Forest Regressor

Ceteris Paribus Profiles



Row 231: A Gentoo penguin

Linear Regression

Ceteris Paribus profile

Advantages

- Works on a **local** level for an individual prediction
- Displays model's **sensitivity** to a feature
- **Intuitive** to understand
- **(Easy to implement)**



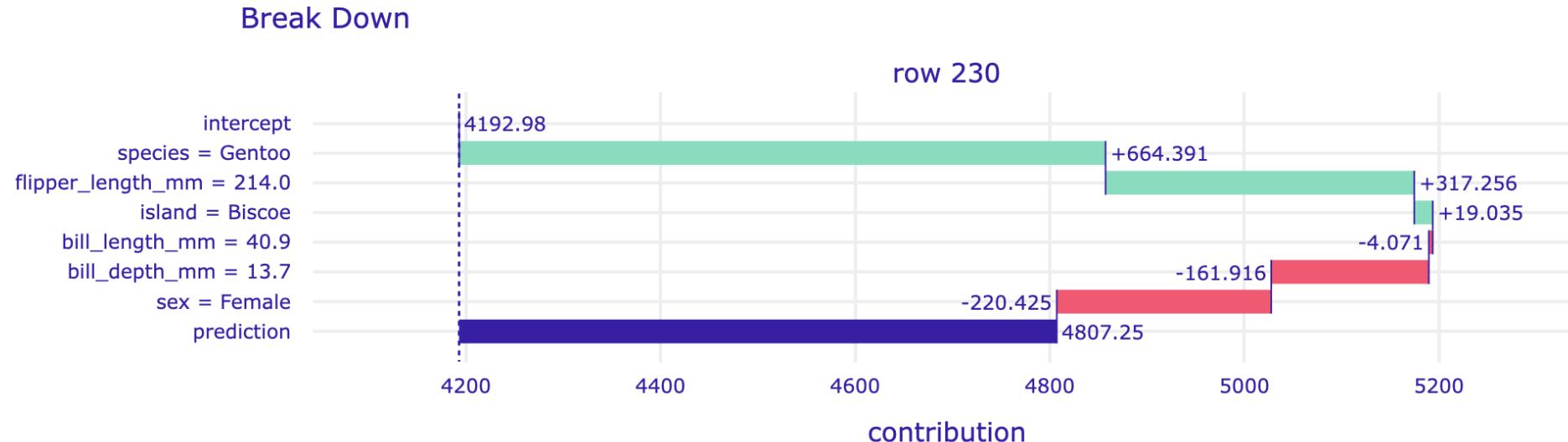
Disadvantages

- Can only display **one feature** meaningfully at the time
- Misleading for **correlated** features
- Only applicable to **one** data point at the time – does not tell you anything about overall trend
- Does not tell you **how much a feature contributes** to individual prediction

Jupyter Notebook 3

03_Prediction_Break_Down.ipynb

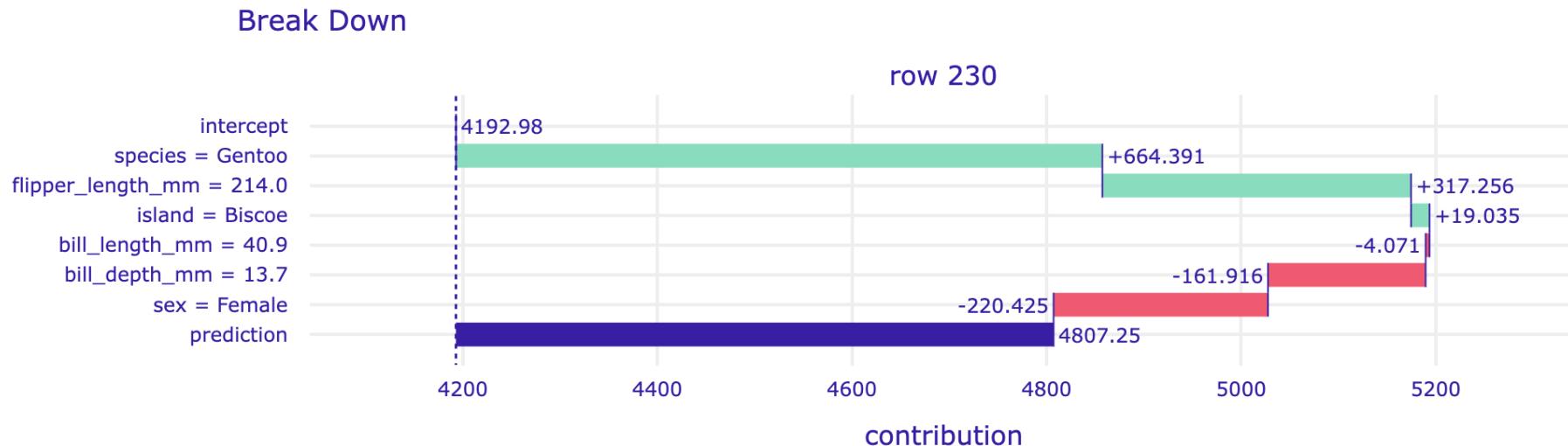
Break-down Plots for Additive Attributions



How to calculate the intercept:

1. Generate predictions with the model on the train set `y_hat = model.predict(X_train)`
2. Calculate the average prediction on the train set `intercept = y_hat.mean()`

Break-down Plots for Additive Attributions



How to calculate the value for *species = Gentoo*:

1. Assume that all datapoints in your dataset `X_train['species'] = 'Gentoo'` have the value Gentoo
2. Calculate the average prediction on the altered (!) train set `value = model.predict(X_train).mean()`
3. Subtract the intercept `value - intercept`

Ceteris Paribus profile

Advantages

- Works on a **local** level for an individual prediction
- Displays a feature's **importance**
- **Fast** to calculate
- **Visual**



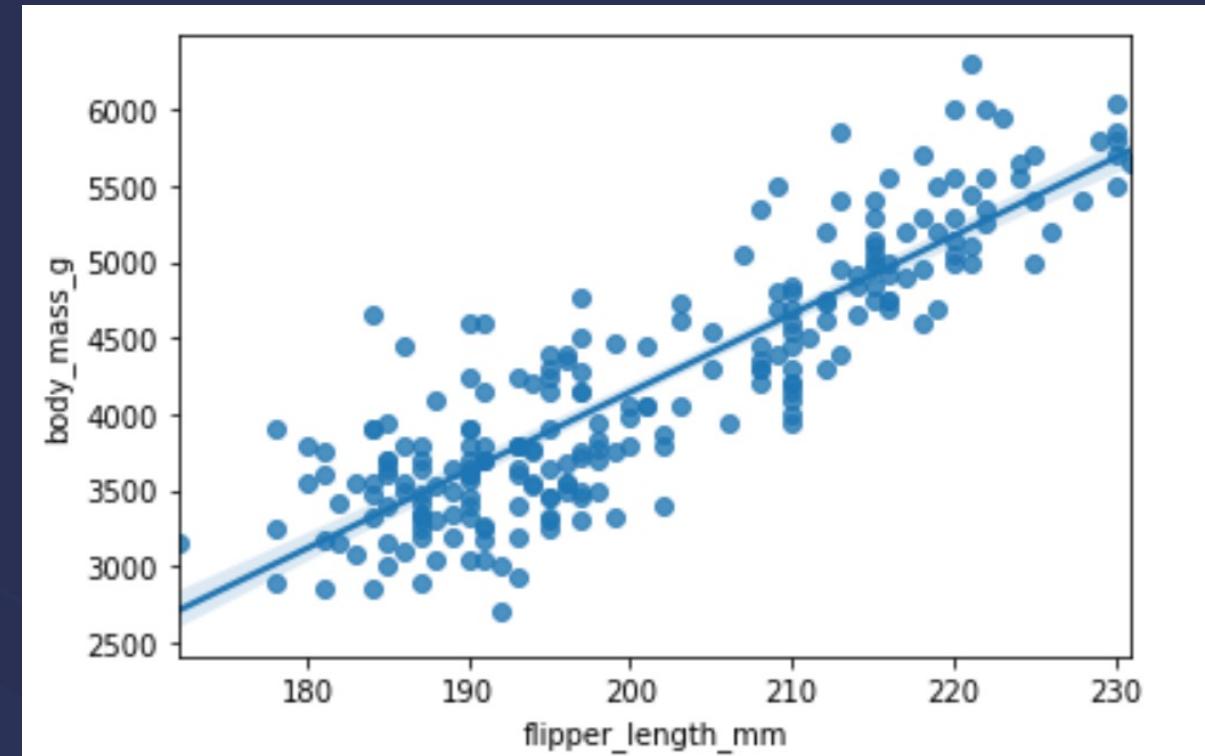
Disadvantages

- The **order** of the features matters
- Misleading for **interactions**
- Messy for many variables
- Does not tell you how **sensitive** a model is to a feature
- Only applicable to **one** data point at the time – does not tell you anything about feature importance

Global Feature Importance

Linear Regression

Global feature
importances



Linear Regression

$m * \text{bill depth} + n * \text{flipper length} + c = \text{body mass}$

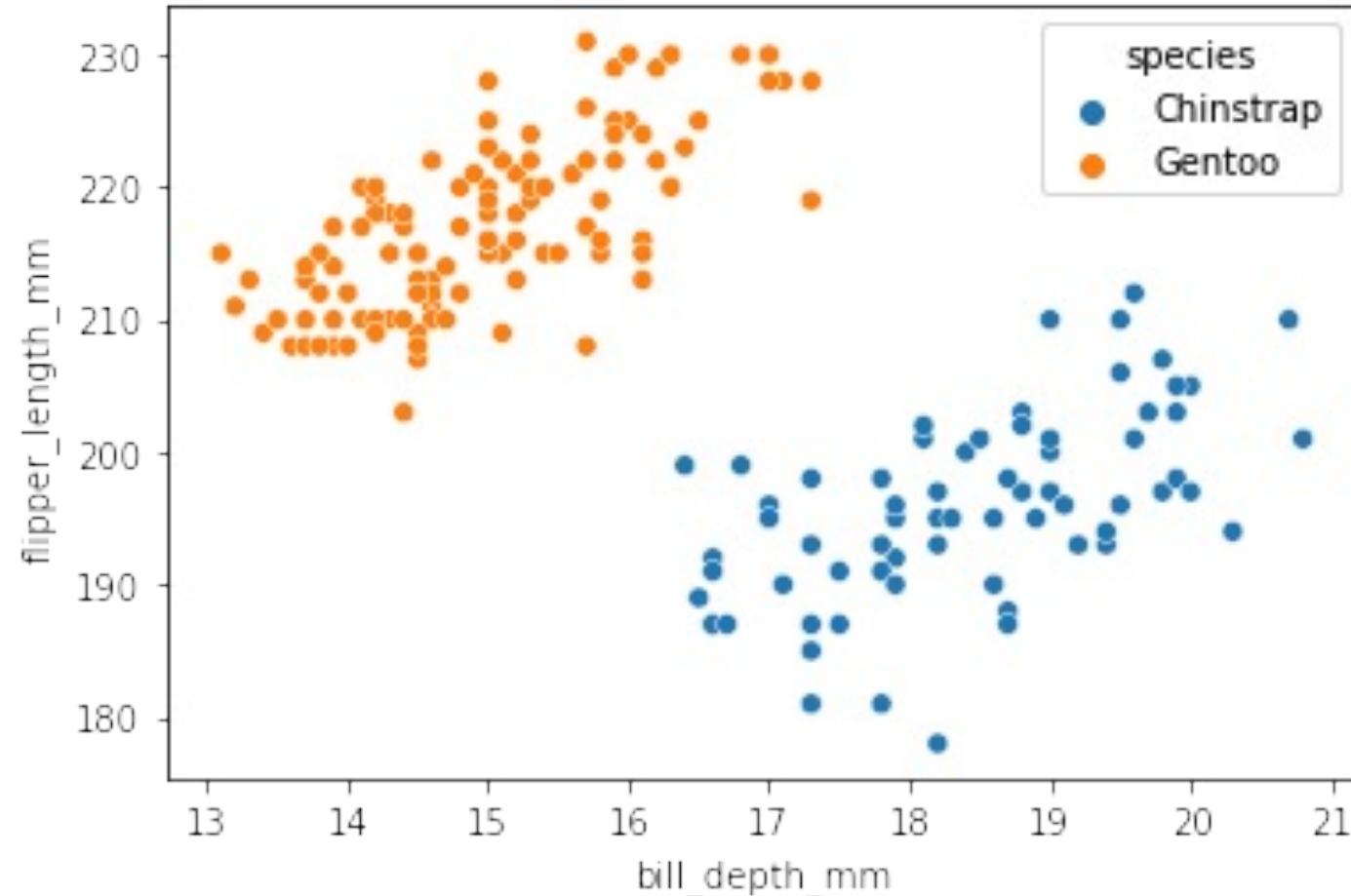
$14 * \text{bill depth} + 53 * \text{flipper length} - 6600 = \text{body mass}$

bill_depth_mm	flipper_length_mm	body_mass_g
18.7	181.0	3750.0

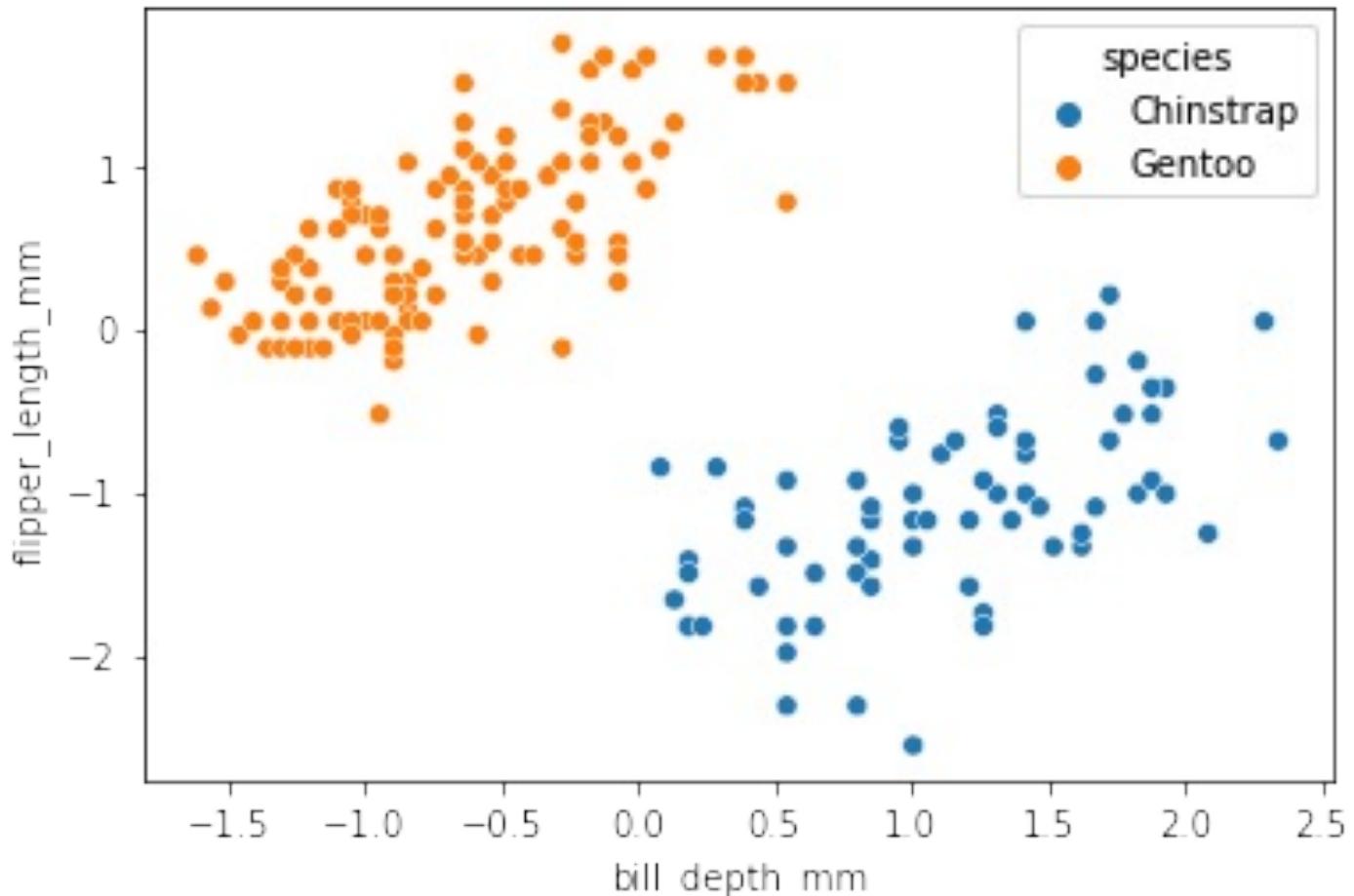
$$14 * 18.7 + 53 * 181.0 - 6600 = 3171$$

Linear Regression

$14 * \text{bill depth} + 53 * \text{flipper length} - 6600 = \text{body mass}$

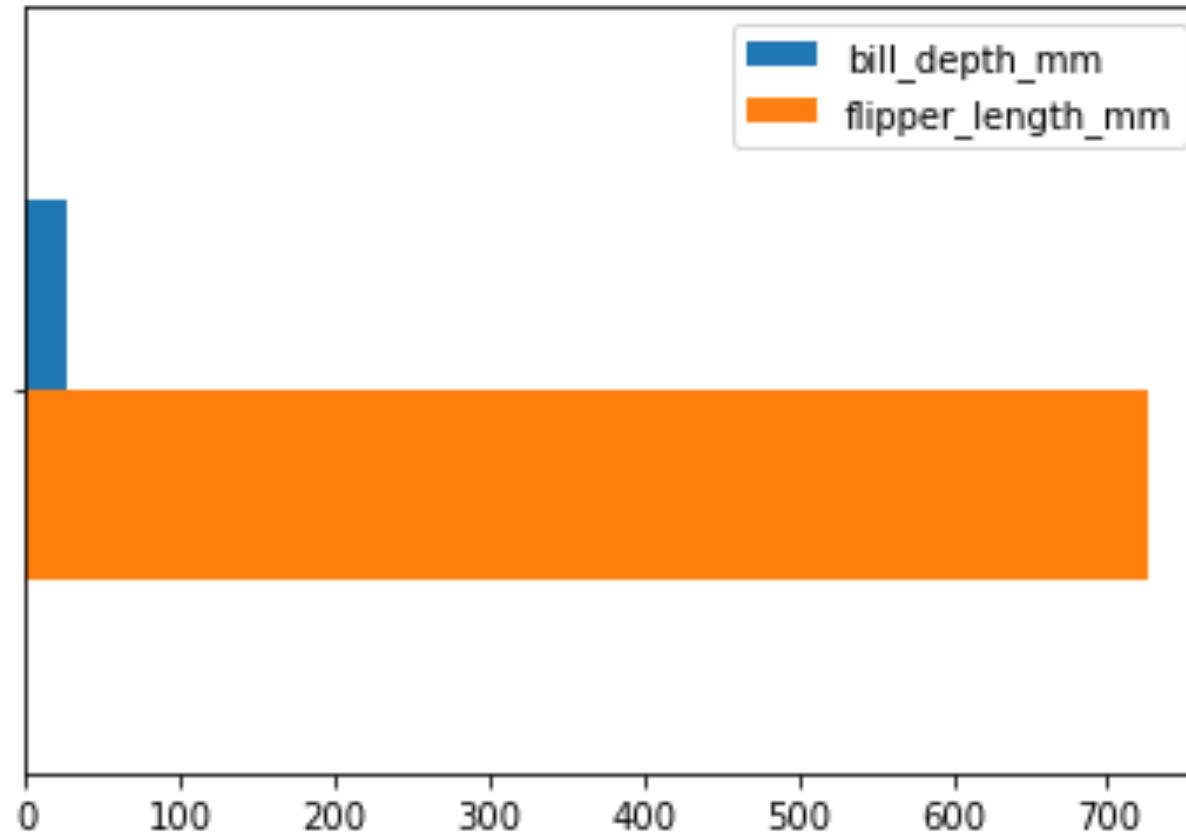


Linear Regression

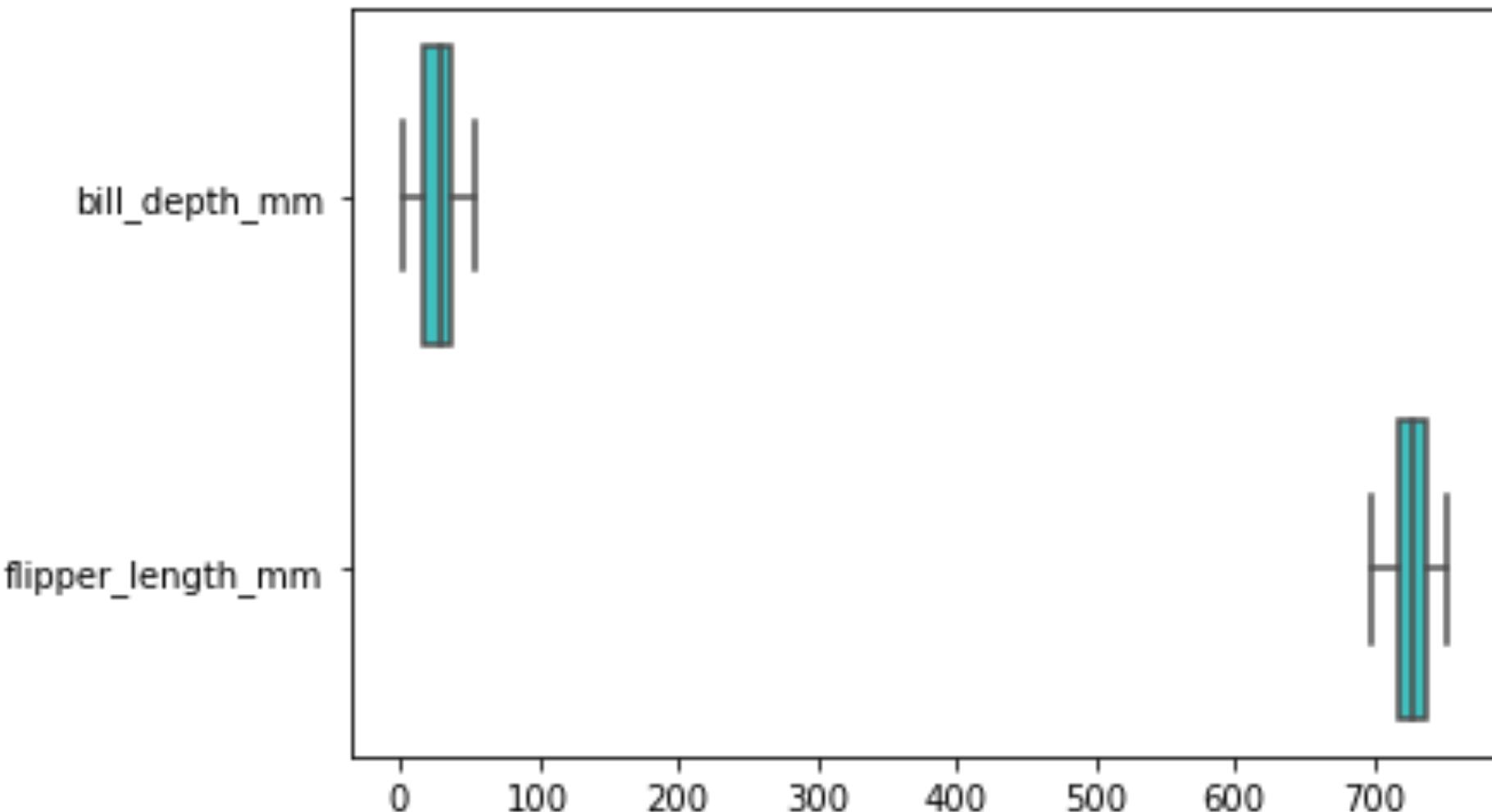


Linear Regression

$28 * \text{bill depth} + 727 * \text{flipper length} + 4199 = \text{body mass}$

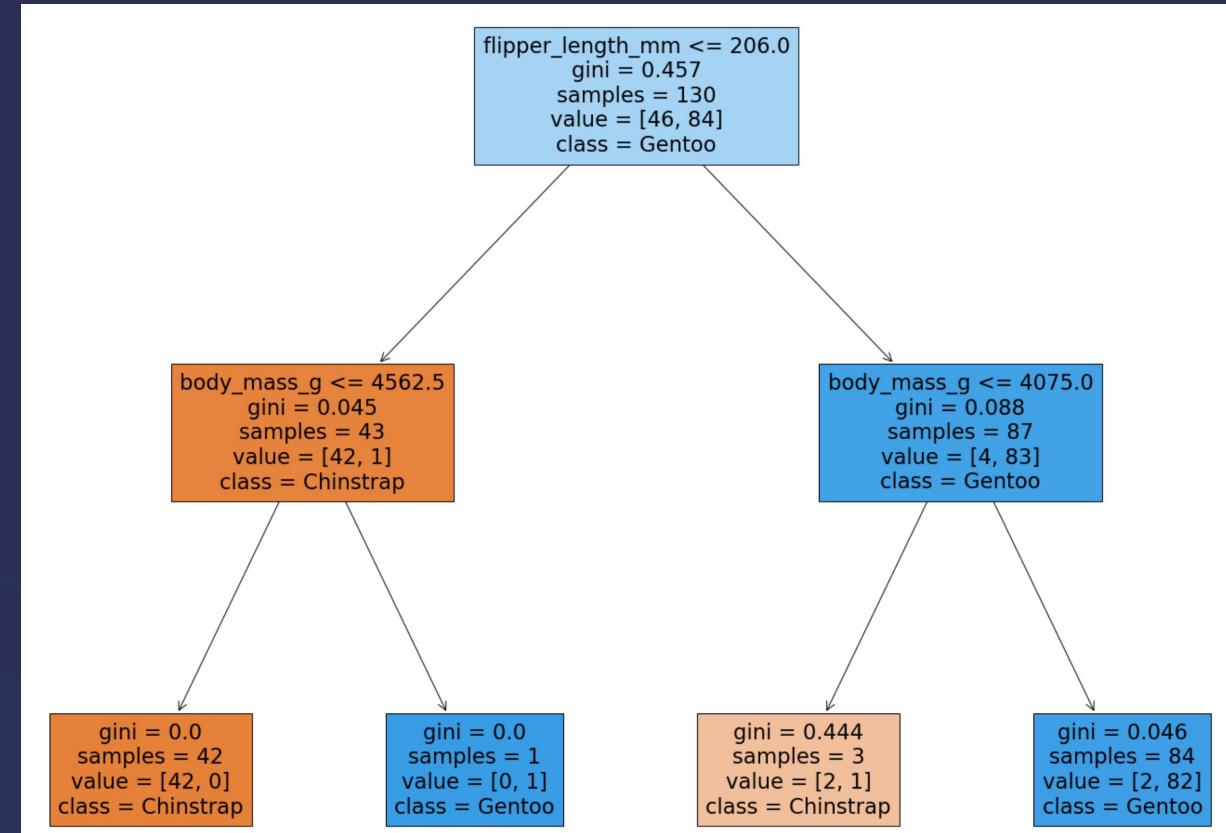


Linear Regression



Decision Tree Classifier

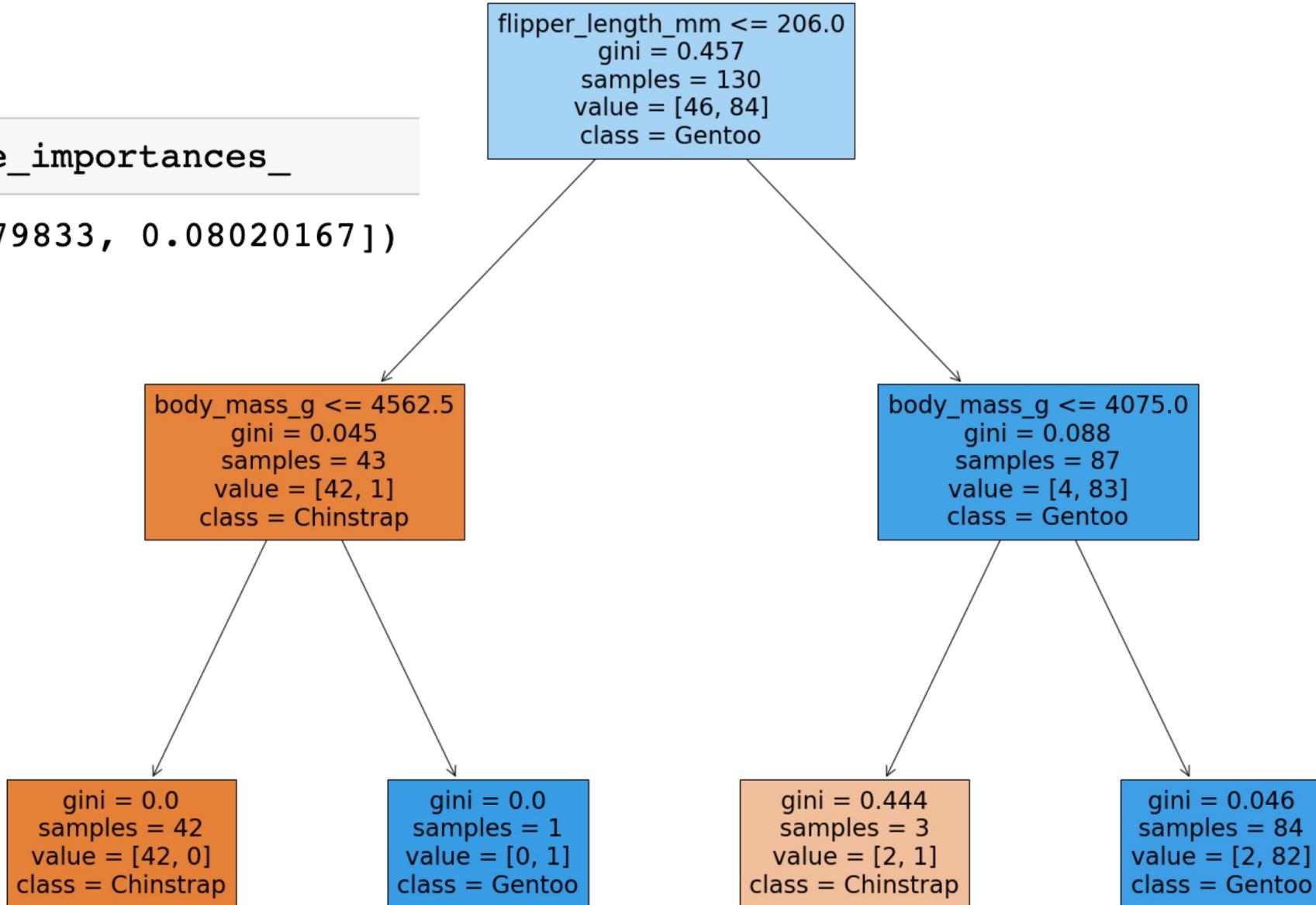
Global feature importances



Decision Tree Classifier

```
model.feature_importances_
```

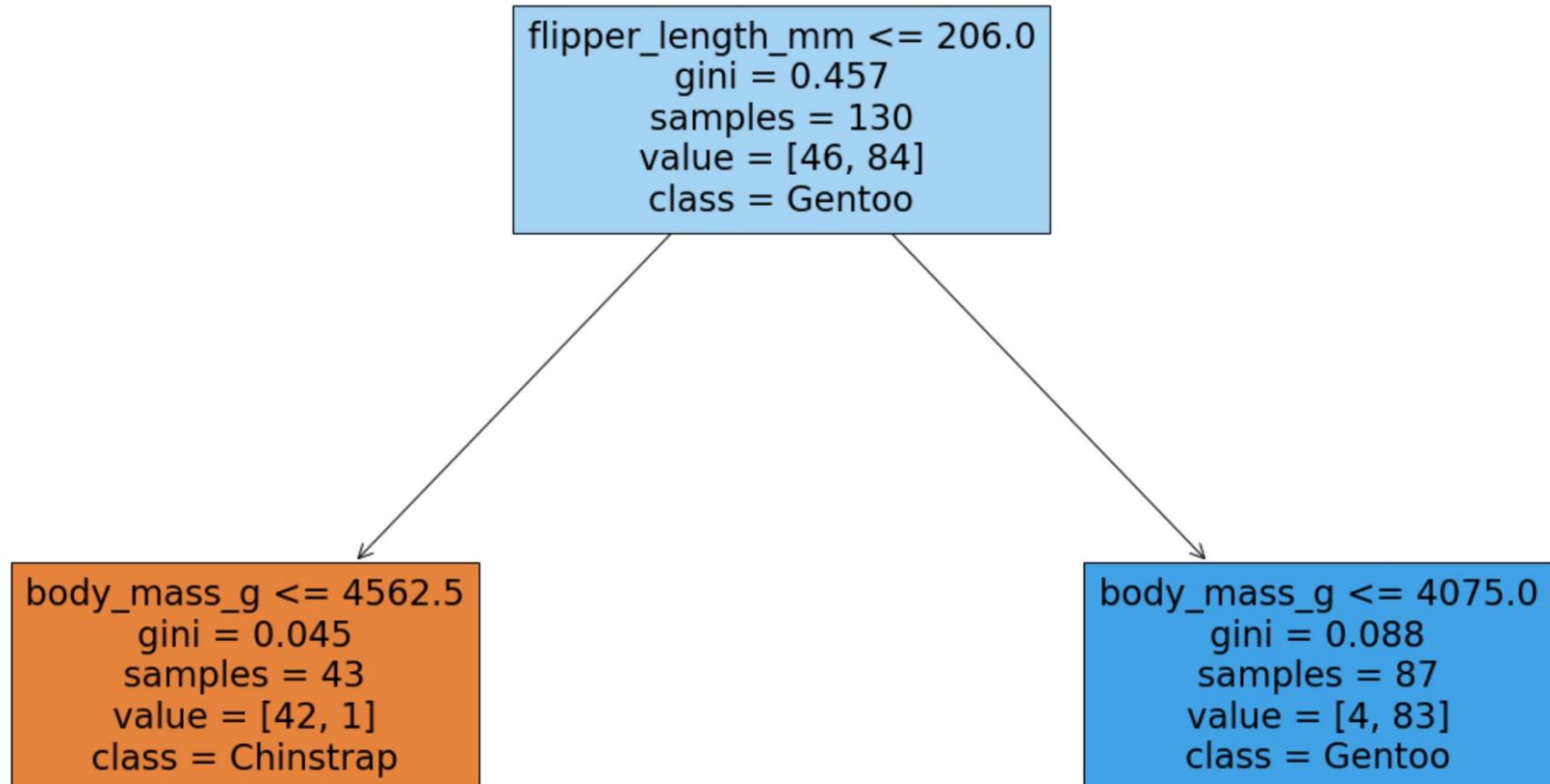
```
array([0.91979833, 0.08020167])
```



Decision Tree Classifier

feature importances

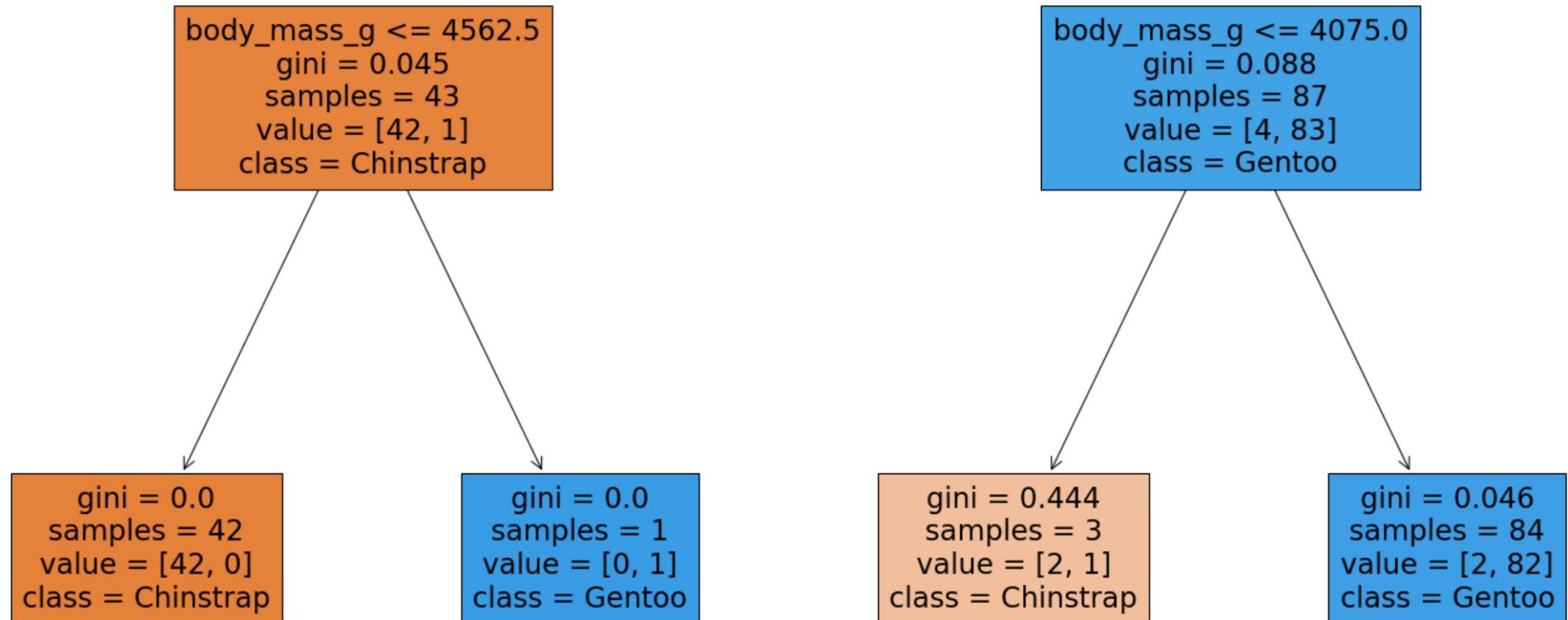
flipper_length_mm	0.919798
body_mass_g	0.080202



Decision Tree Classifier

feature importances

flipper_length_mm	0.919798
body_mass_g	0.080202



Random Forest Classifier

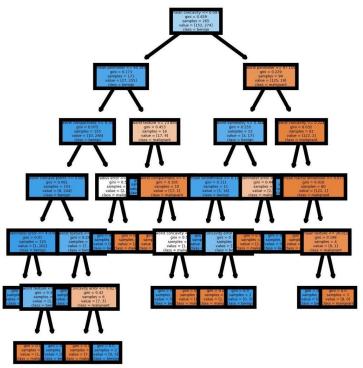
GO
DATA
DRIVEN

**Global feature
importances**

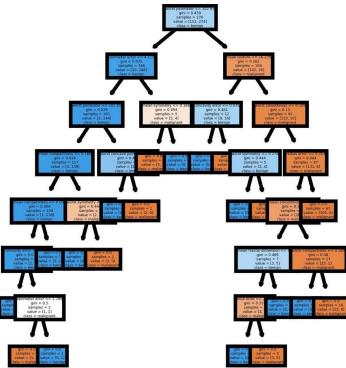


Random Forest Classifier

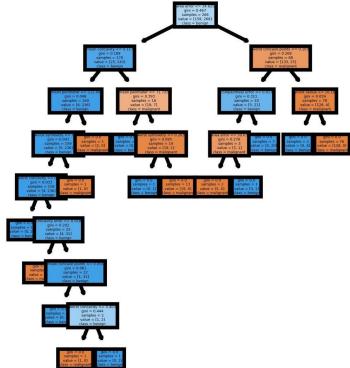
Estimator: 0



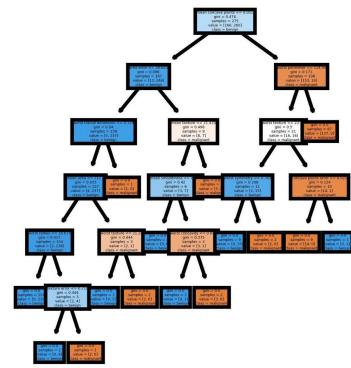
Estimator: 1



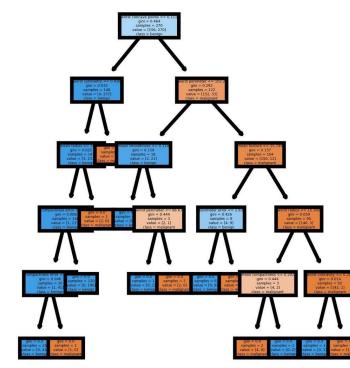
Estimator: 2



Estimator: 3



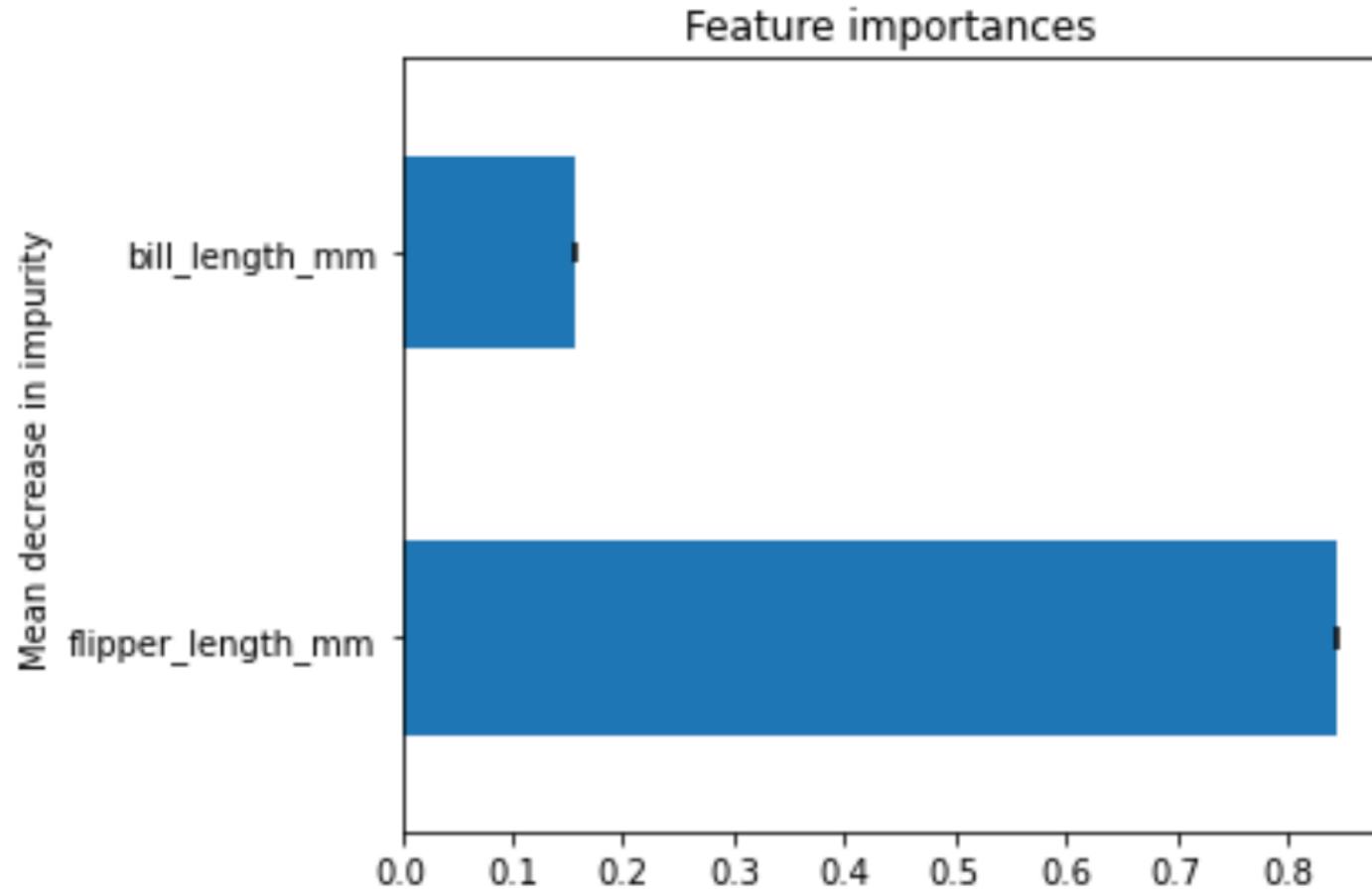
Estimator: 4



```
model.feature_importances_
```

```
array([ 0.84263378,  0.15736622])
```

Random Forest Classifier



Jupyter Notebook 4

04_Permutation_Feature_Importance.
ipy_nb

GO 
DATA
DRIVEN

Model-Agnostic Feature Importance

	flipper_length_mm	bill_length_mm	bill_depth_mm	sex	island	species
27	187.0	40.5	17.9	Female	Biscoe	Adelie
292	221.0	48.2	15.6	Male	Biscoe	Gentoo
302	212.0	47.4	14.6	Female	Biscoe	Gentoo
62	185.0	37.6	17.0	Female	Biscoe	Adelie
181	205.0	52.8	20.0	Male	Dream	Chinstrap
15	185.0	36.6	17.8	Female	Torgersen	Adelie
231	216.0	49.0	16.1	Male	Biscoe	Gentoo
233	213.0	48.4	14.6	Male	Biscoe	Gentoo
87	189.0	36.9	18.6	Female	Dream	Adelie
159	197.0	51.3	18.2	Male	Dream	Chinstrap



0.856

Drop Feature Importance

	flipper_length_mm	bill_length_mm	bill_depth_mm	sex	island	species
27	187.0	40.5	17.9	Female	Biscoe	Adelie
292	221.0	48.2	15.6	Male	Biscoe	Gentoo
302	212.0	47.4	14.6	Female	Biscoe	Gentoo
62	185.0	37.6	17.0	Female	Biscoe	Adelie
181	205.0	52.8	20.0	Male	Dream	Chinstrap
15	185.0	36.6	17.8	Female	Torgersen	Adelie
231	216.0	49.0	16.1	Male	Biscoe	Gentoo
233	213.0	48.4	14.6	Male	Biscoe	Gentoo
87	189.0	36.9	18.6	Female	Dream	Adelie
159	197.0	51.3	18.2	Male	Dream	Chinstrap



	flipper_length_mm	bill_length_mm	sex	island	species
27	187.0	40.5	Female	Biscoe	Adelie
292	221.0	48.2	Male	Biscoe	Gentoo
302	212.0	47.4	Female	Biscoe	Gentoo
62	185.0	37.6	Female	Biscoe	Adelie
181	205.0	52.8	Male	Dream	Chinstrap
15	185.0	36.6	Female	Torgersen	Adelie
231	216.0	49.0	Male	Biscoe	Gentoo
233	213.0	48.4	Male	Biscoe	Gentoo
87	189.0	36.9	Female	Dream	Adelie
159	197.0	51.3	Male	Dream	Chinstrap

0.856

???

Permutation Feature Importance

	flipper_length_mm	bill_length_mm	bill_depth_mm	sex	island	species
27	187.0	40.5	17.9	Female	Biscoe	Adelie
292	221.0	48.2	15.6	Male	Biscoe	Gentoo
302	212.0	47.4	14.6	Female	Biscoe	Gentoo
62	185.0	37.6	17.0	Female	Biscoe	Adelie
181	205.0	52.8	20.0	Male	Dream	Chinstrap
15	185.0	36.6	17.8	Female	Torgersen	Adelie
231	216.0	49.0	16.1	Male	Biscoe	Gentoo
233	213.0	48.4	14.6	Male	Biscoe	Gentoo
87	189.0	36.9	18.6	Female	Dream	Adelie
159	197.0	51.3	18.2	Male	Dream	Chinstrap



	flipper_length_mm	bill_length_mm	bill_depth_mm	sex	island	species
27	187.0	40.5	17.9	Female	Biscoe	Adelie
292	221.0	48.2	15.6	Male	Biscoe	Gentoo
302	212.0	47.4	14.6	Female	Biscoe	Gentoo
62	185.0	37.6	17.0	Female	Biscoe	Adelie
181	205.0	52.8	20.0	Male	Dream	Chinstrap
15	185.0	36.6	17.8	Female	Torgersen	Adelie
231	216.0	49.0	16.1	Male	Biscoe	Gentoo
233	213.0	48.4	14.6	Male	Biscoe	Gentoo
87	189.0	36.9	18.6	Female	Dream	Adelie
159	197.0	51.3	18.2	Male	Dream	Chinstrap



0.856

Permutation Feature Importance

Advantages

- Works on a **global** level for a model
- Displays a feature's **importance**
- **Fast** to calculate
- **Additional** tricks



Disadvantages

- Dependent on **randomness**
- Misleading for **correlations**
- Can lead to **unrealistic** data points
- Need access to **true** outcome
- Does not tell you how **sensitive** a model is to a feature

Jupyter Notebook 5

05_Partial_Dependence_Plots.ipynb

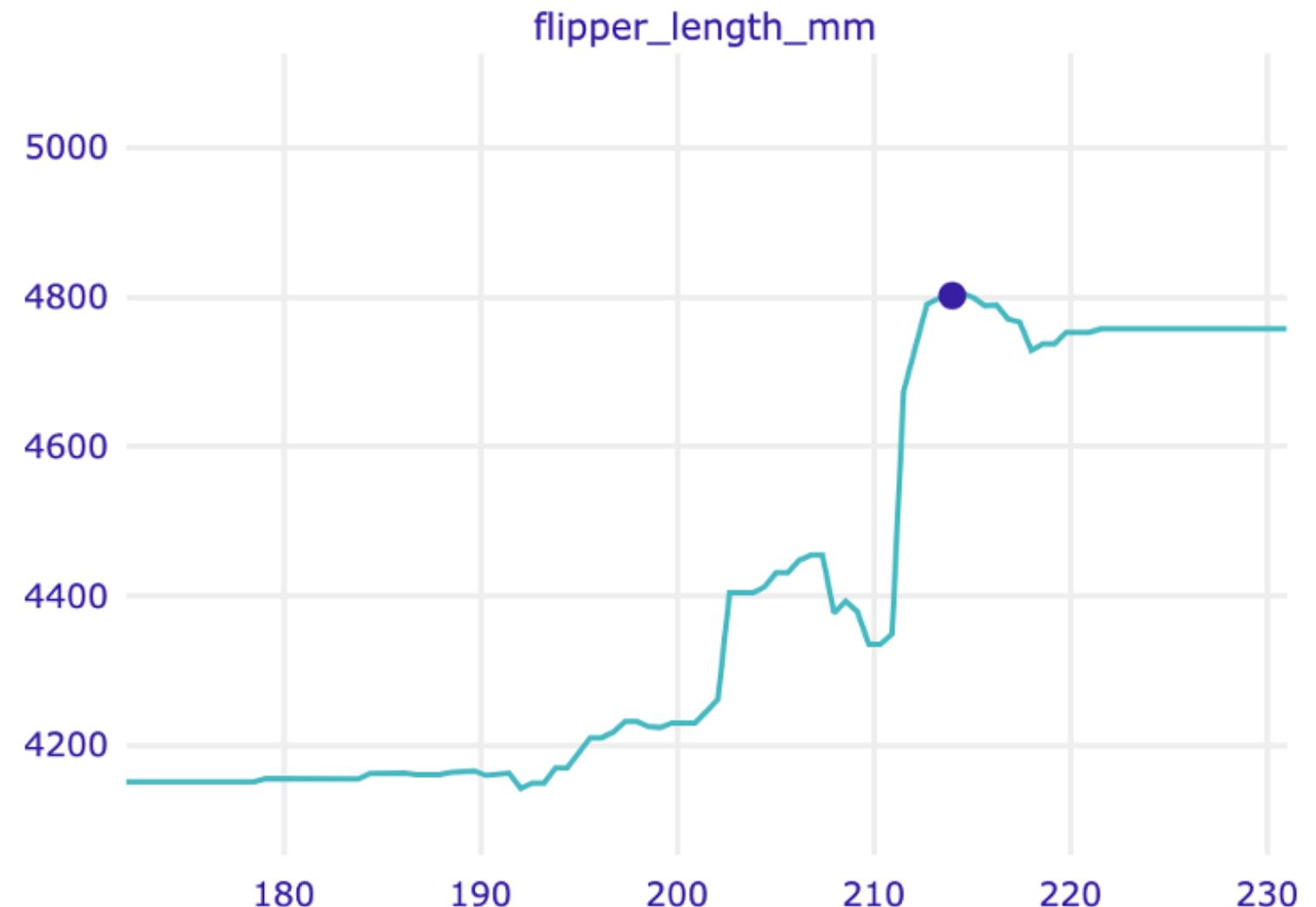
GO 
DATA
DRIVEN

Global Sensitivity

Ceteris Paribus

- Sensivity to feature
- On a local level

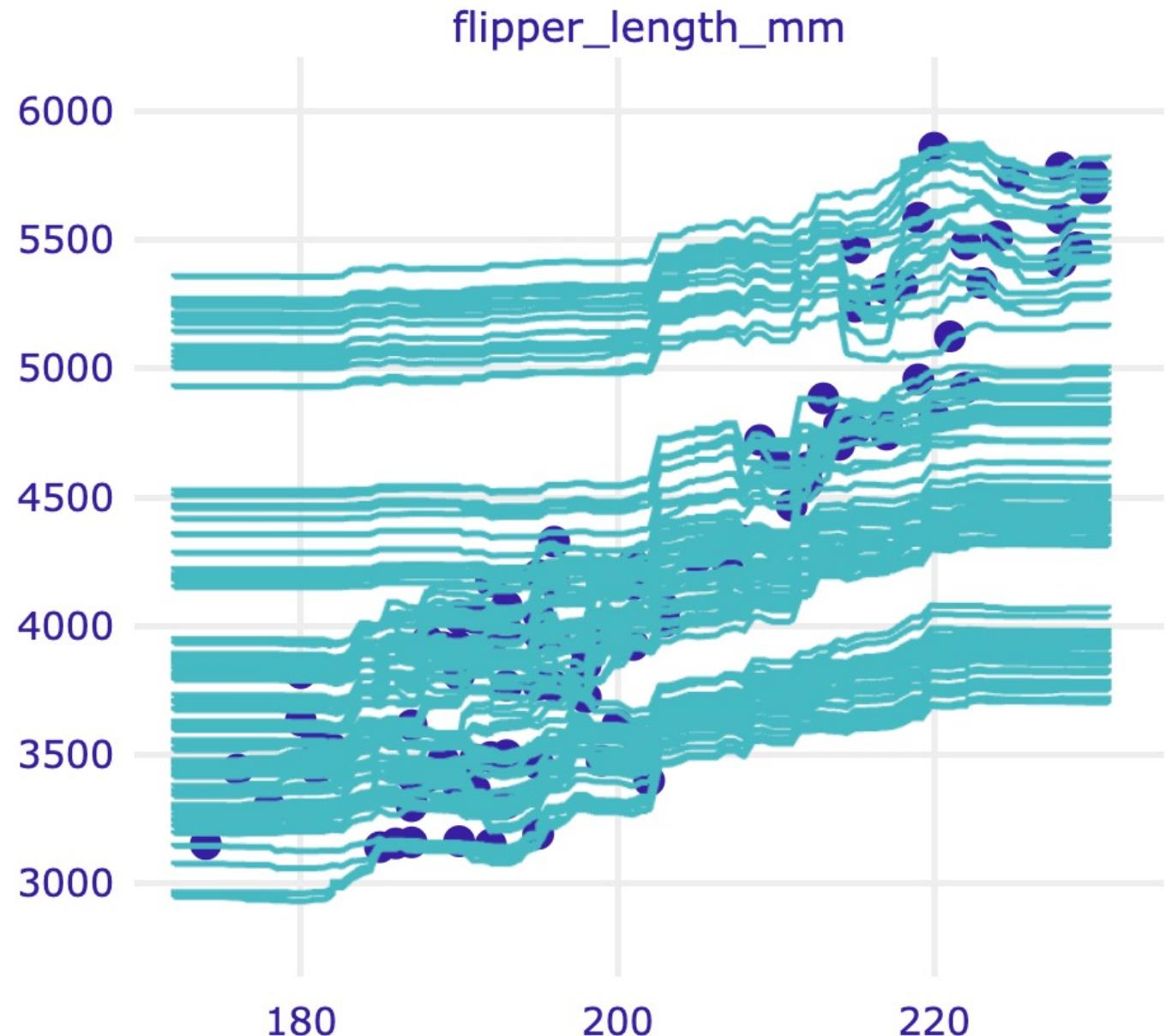
Ceteris Paribus Profiles



Individual Conditional Expectation

- Sensivity to feature
- One line per instance

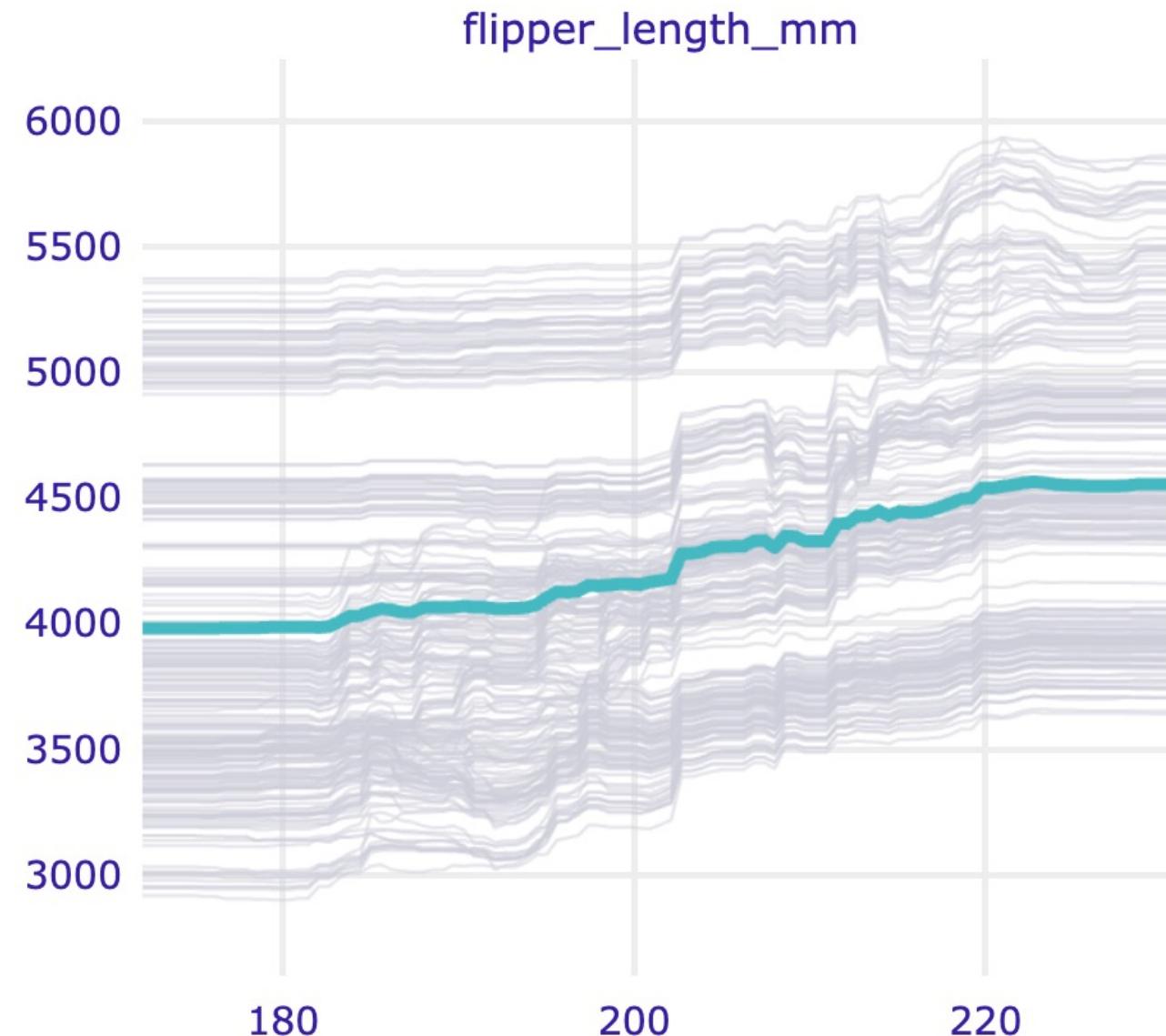
Ceteris Paribus Profiles



Partial Dependence Plot

- Sensivity to feature
- Average of ICE plot
- Global method

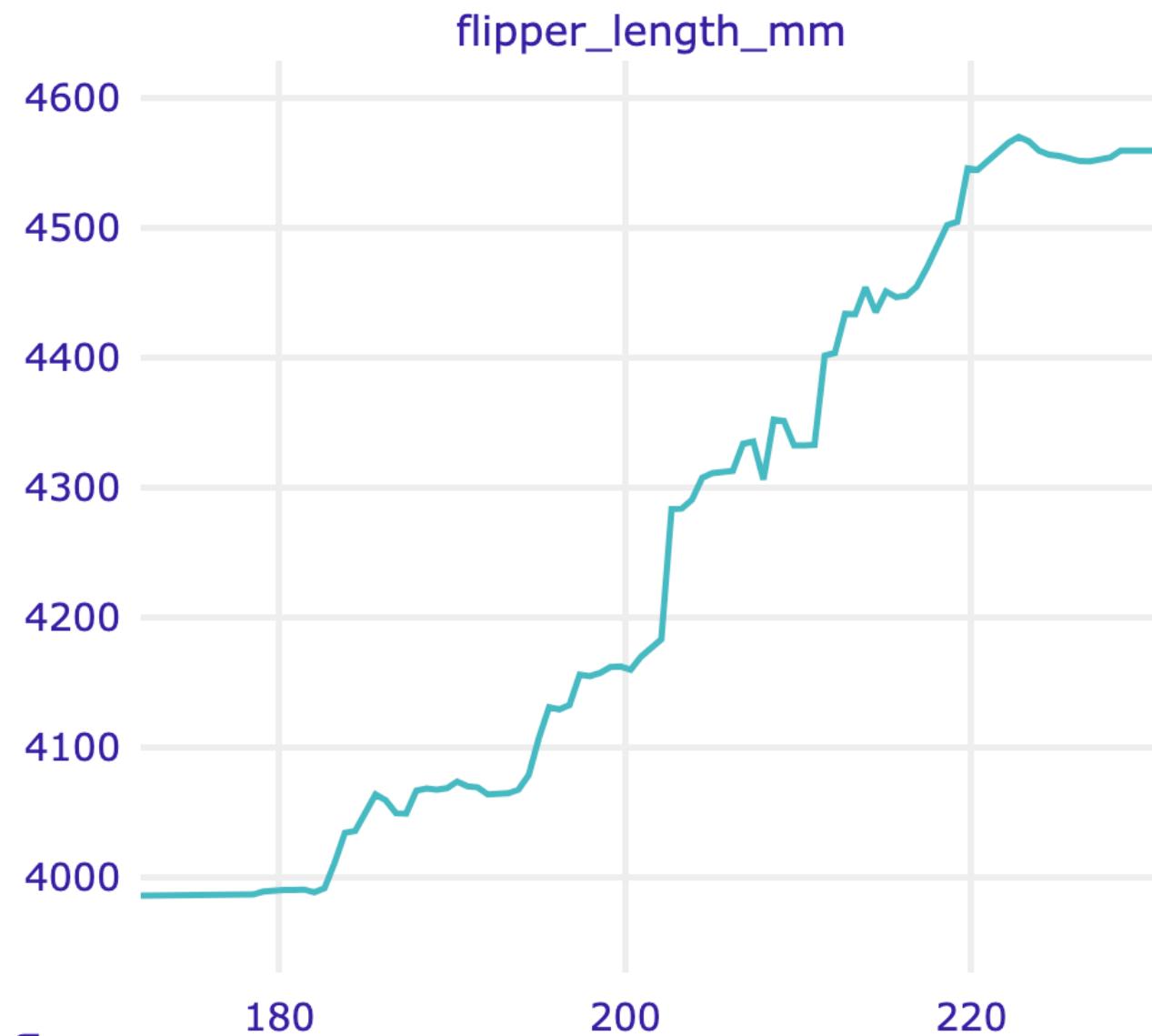
Aggregated Profiles



Partial Dependence Plot

- Sensivity to feature
- Average of ICE plot
- Global method

Aggregated Profiles



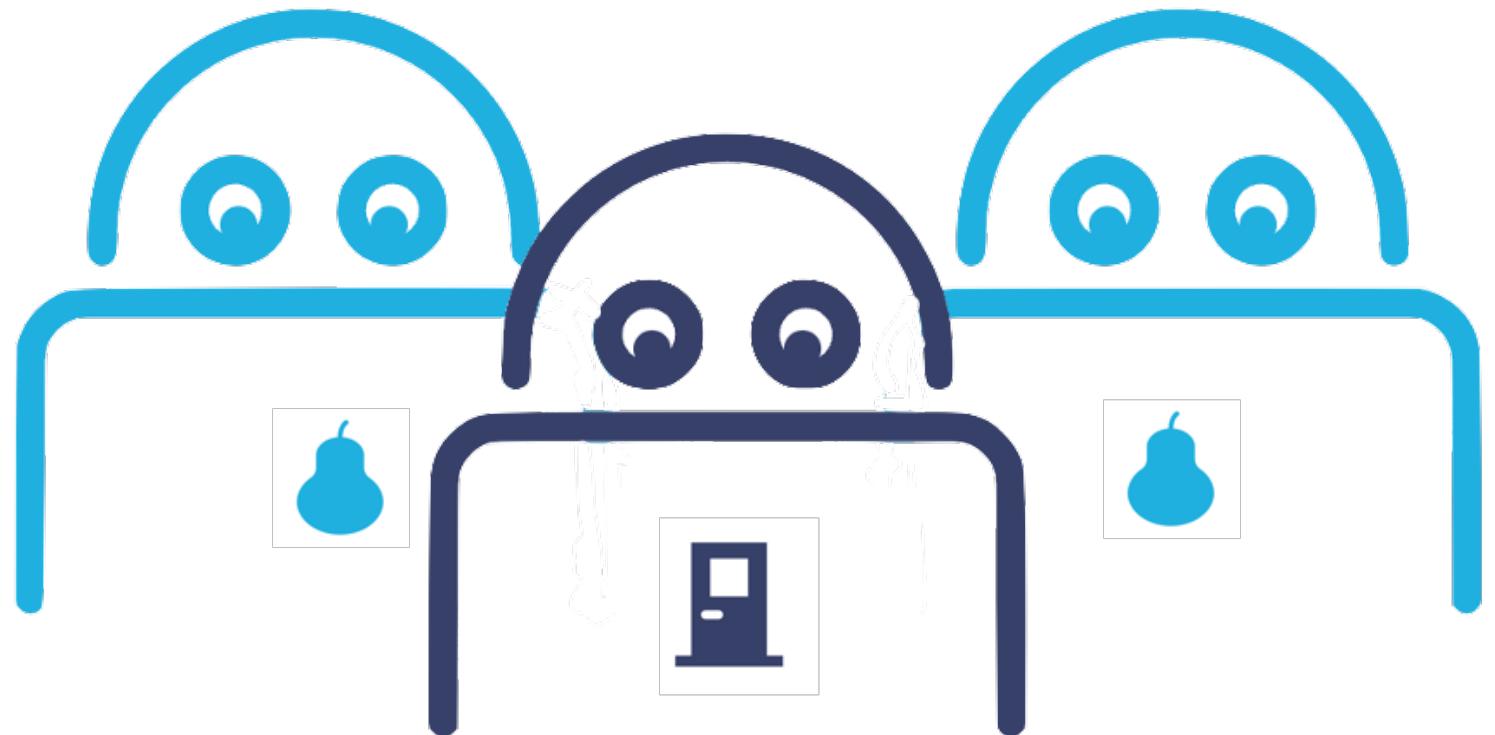
Summary

- Intrinsicly interpretable machine learning models
- Model-agnostic interpretability techniques

	Local	Global
Importance	Break-down plots	Permutation Feature Importance
Sensitivity	Ceteris Paribus	Partial Dependence Plots

Hackathon!

- Credit Lending
- Laptop Prices
- Medical Diagnoses
- Discussion Hackathon
(theory)



GO 
DATA
DRIVEN

Unpacking the Black Box

How to Interpret your Machine Learning Model

GoDataDriven

@ AMLD Switzerland 2022

GO
DATA
DRIVEN

The Data & AI Maturity Journey
13:30-17:30, March 28 @ 1A

GoDataDriven Academy

- Public trainings
- In-company tailored trainings
- Virtual & in-person

 Visit godatadriven.academy

GO 
DATA
DRIVEN



GO 
DATA
DRIVEN

Unpacking the Black Box

How to Interpret your Machine Learning Model