

Computational Analysis of Big Data

Week 5

Machine Learning 2

Decision trees

Ensemble learning

Logistic regression

Recap: Machine Learning

Canonical example



[]

Dog



[fluffy sad looking showing teeth ears down tail between legs
 1 1 0 1 0 0 0 0]
is retriever growling

[]

Data point

Canonical example



[]

Dog



fluffy
sad looking
showing teeth
ears down
tail between legs
is retriever
growling
[0]

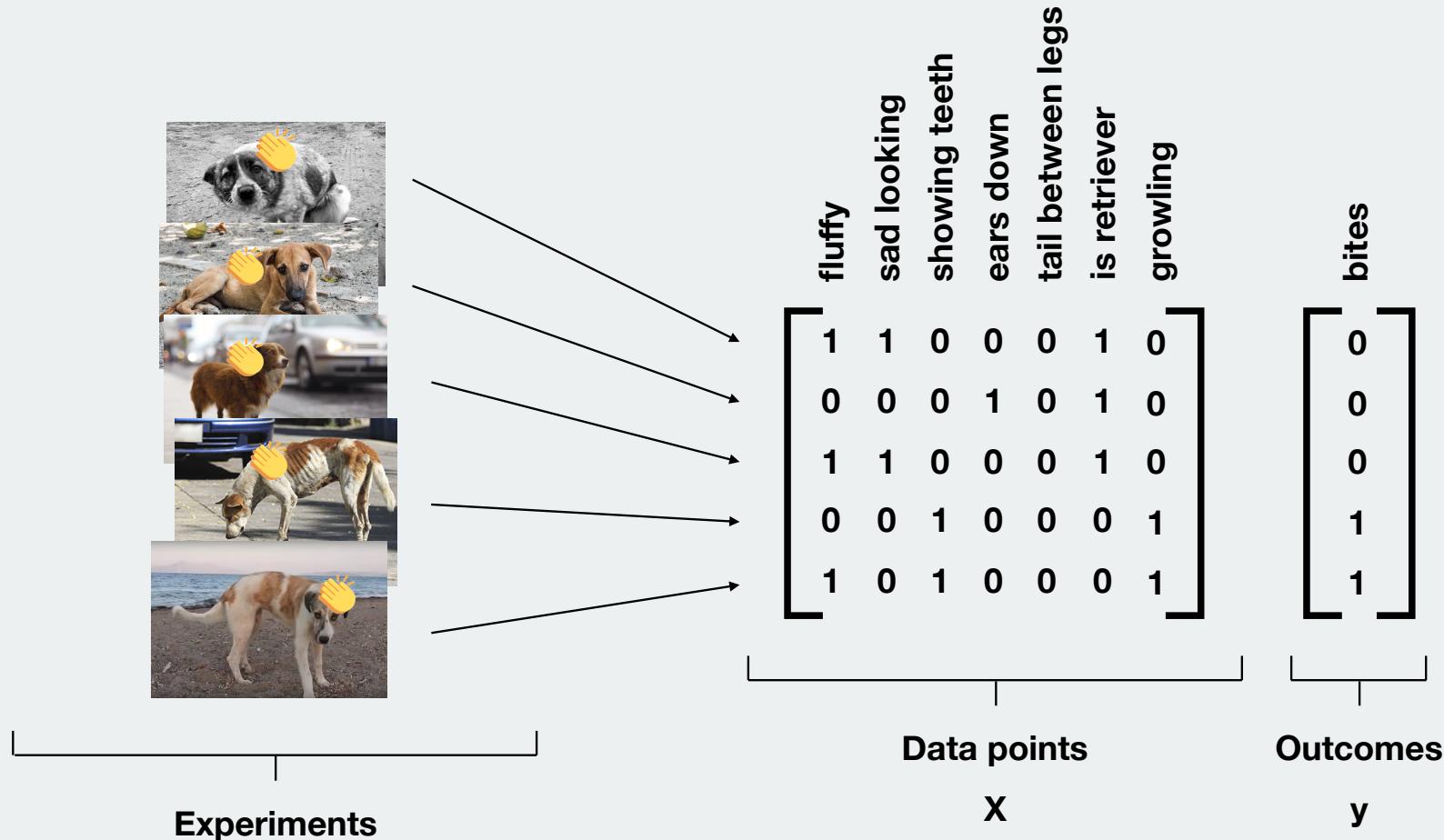
[]

Data point

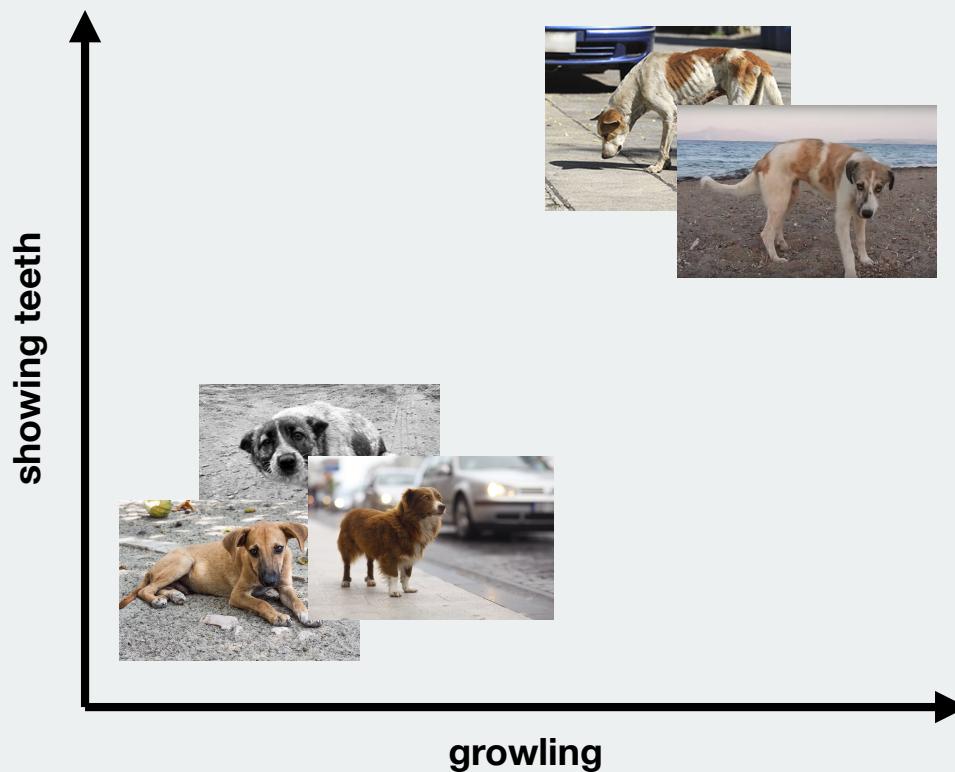
[]

Outcome

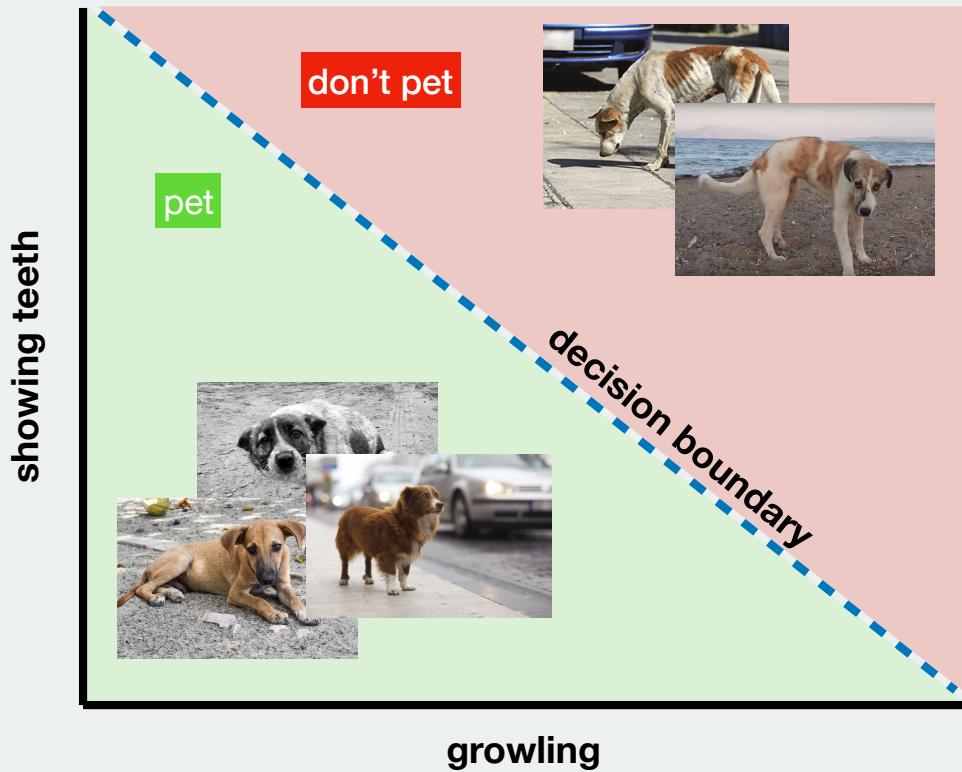
Canonical example



Canonical example



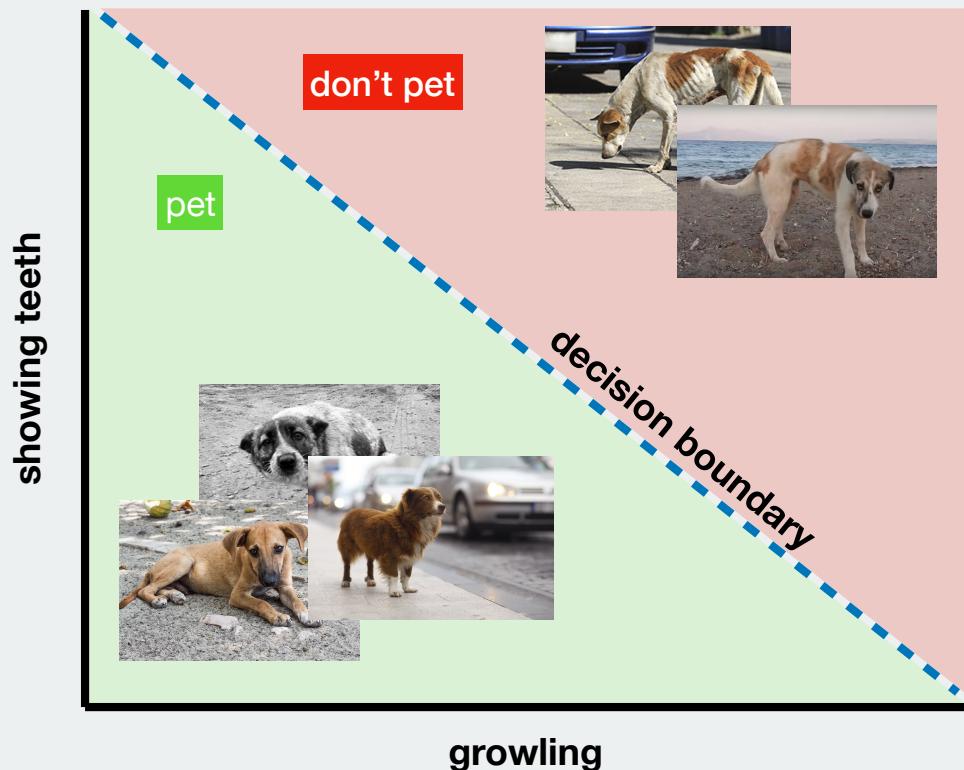
Canonical example



Canonical example

Supervised Machine Learning

When the input data has outcome labels



Decision trees

Ensemble learning

Logistic regression

Decision trees

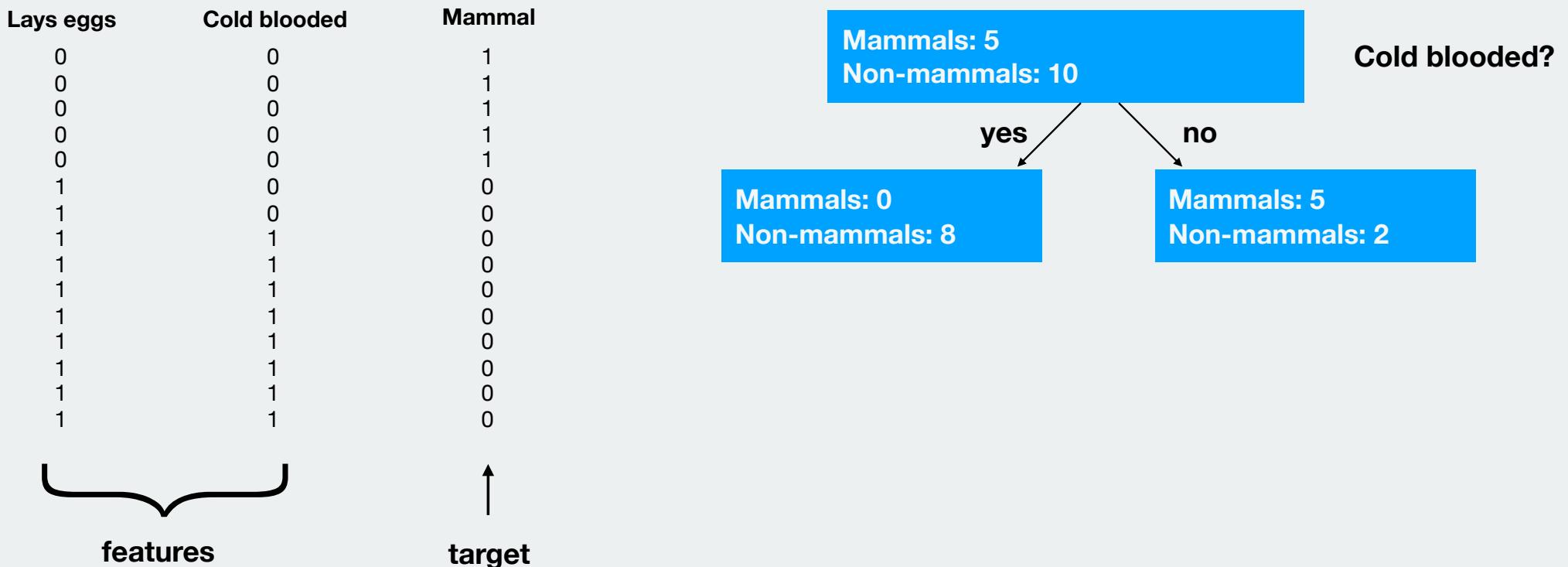
Decision trees

Lays eggs	Cold blooded	Mammal	
0	0	1	
0	0	1	
0	0	1	
0	0	1	
0	0	1	
1	0	0	
1	0	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	
1	1	0	

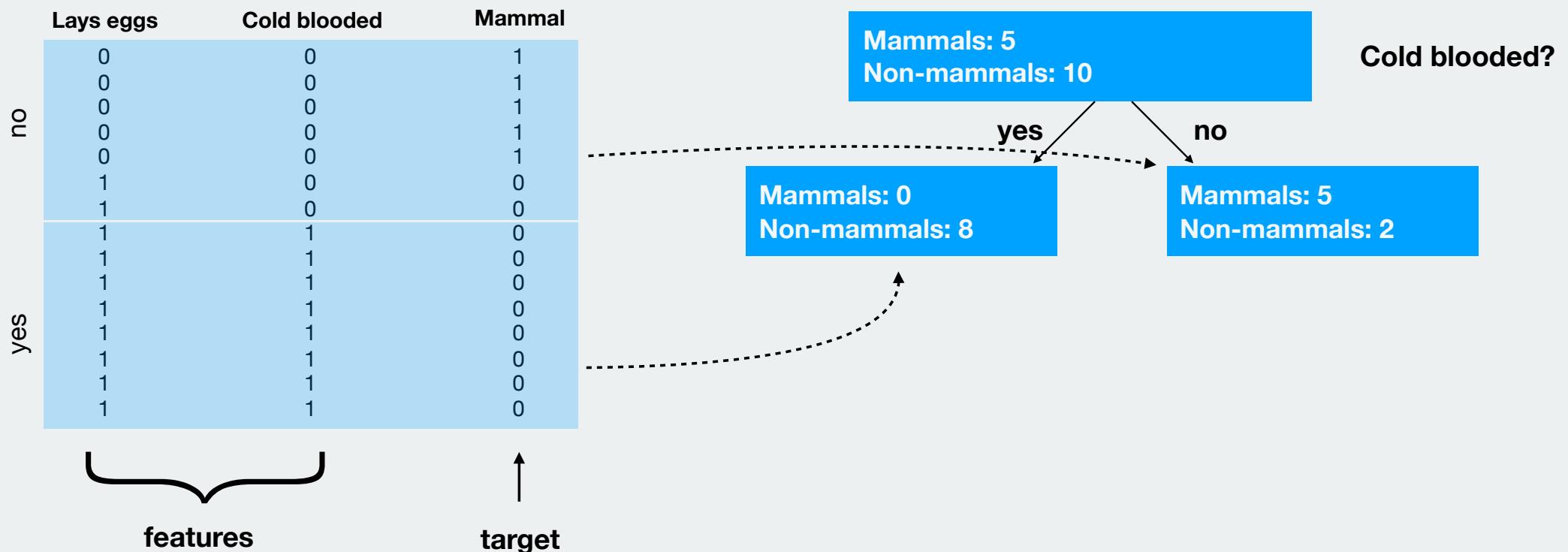
features **target**

Mammals: 5
Non-mammals: 10

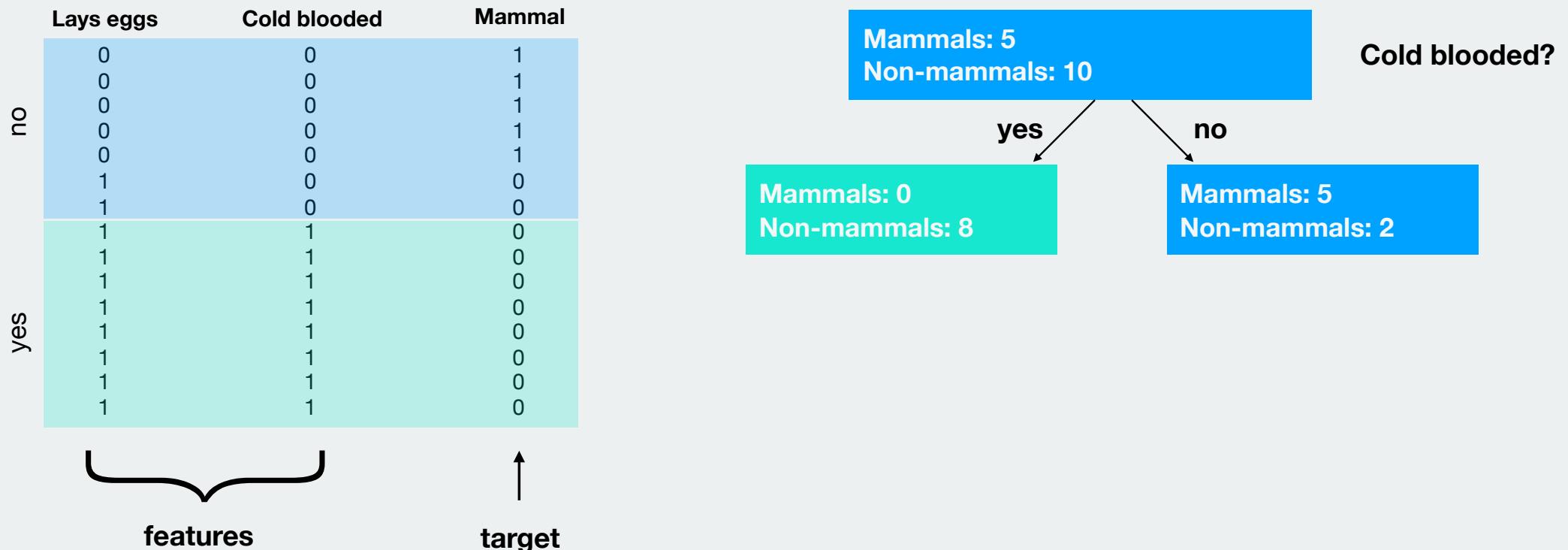
Decision trees



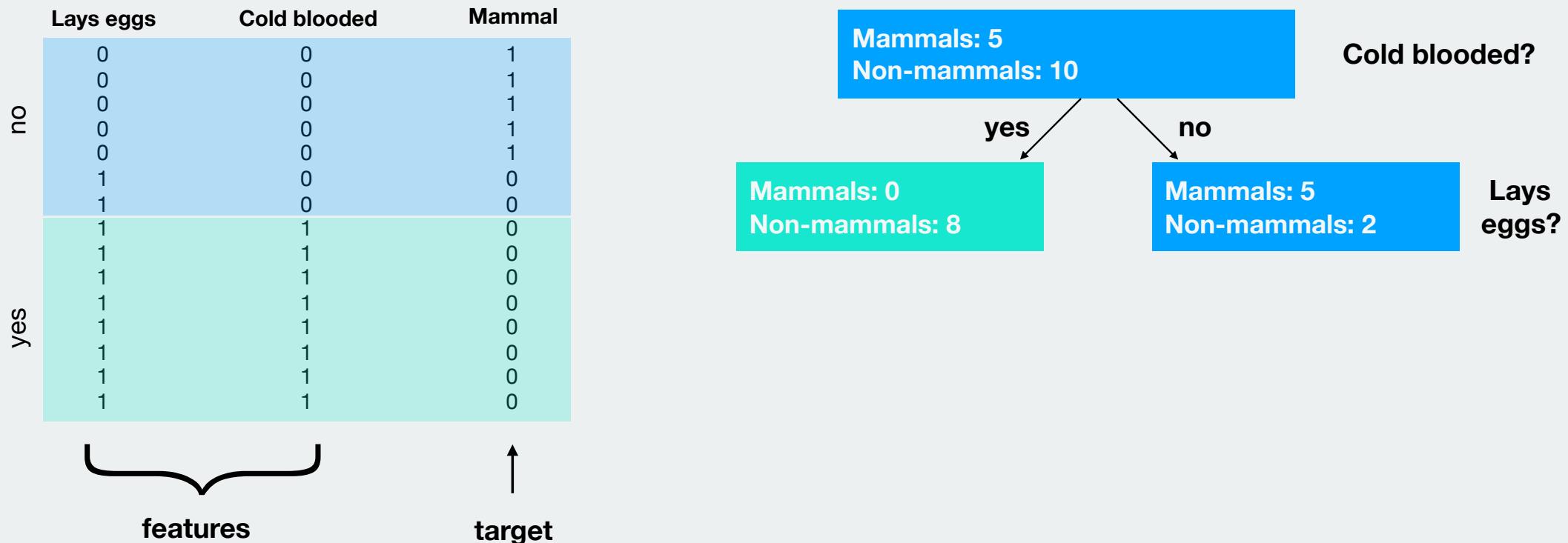
Decision trees



Decision trees



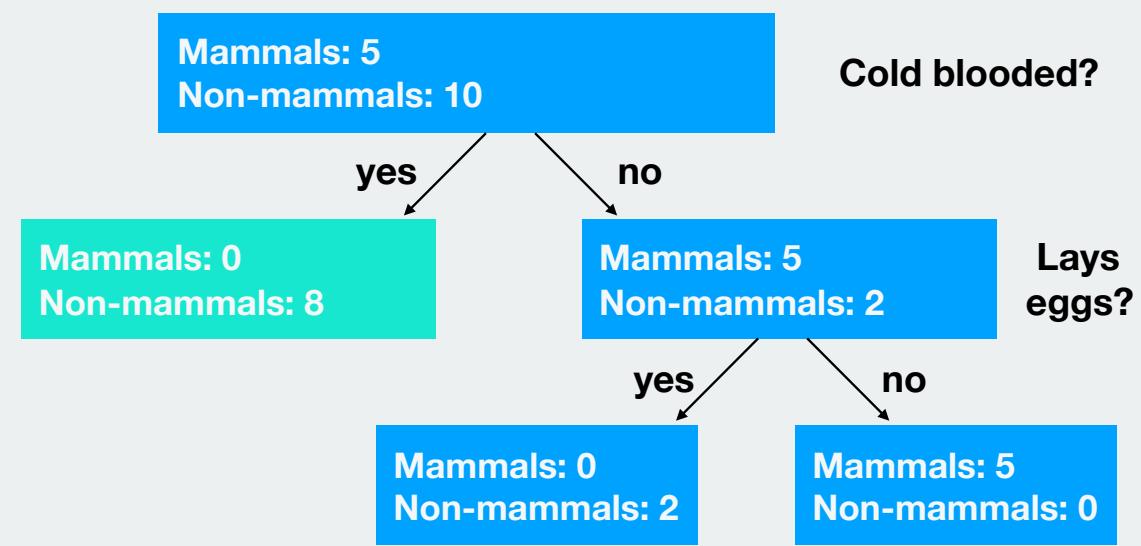
Decision trees



Decision trees

	Lays eggs	Cold blooded	Mammal
no	0 0 0 0 0	0 0 0 0 0	1 1 1 1 1
yes	1 1	0 0	0 0
	1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1	0 0 0 0 0 0 0 0

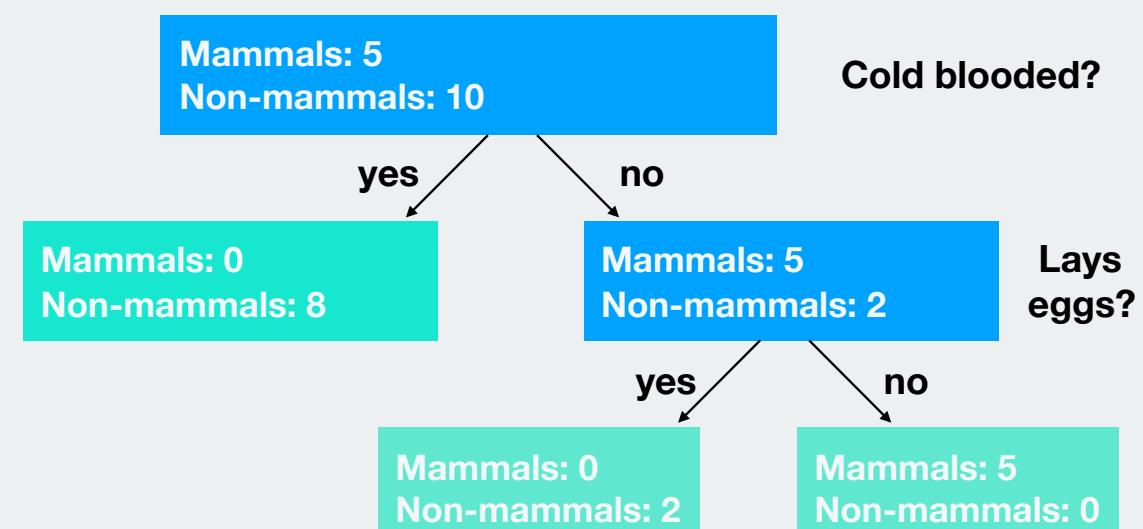
↑
features
target



Decision trees

	Lays eggs	Cold blooded	Mammal
no	0 0 0 0 0 0	0 0 0 0 0 0	1 1 1 1 1 1
yes	1 1	0 0	0 0
	1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1	0 0 0 0 0 0 0 0 0 0

features **target**



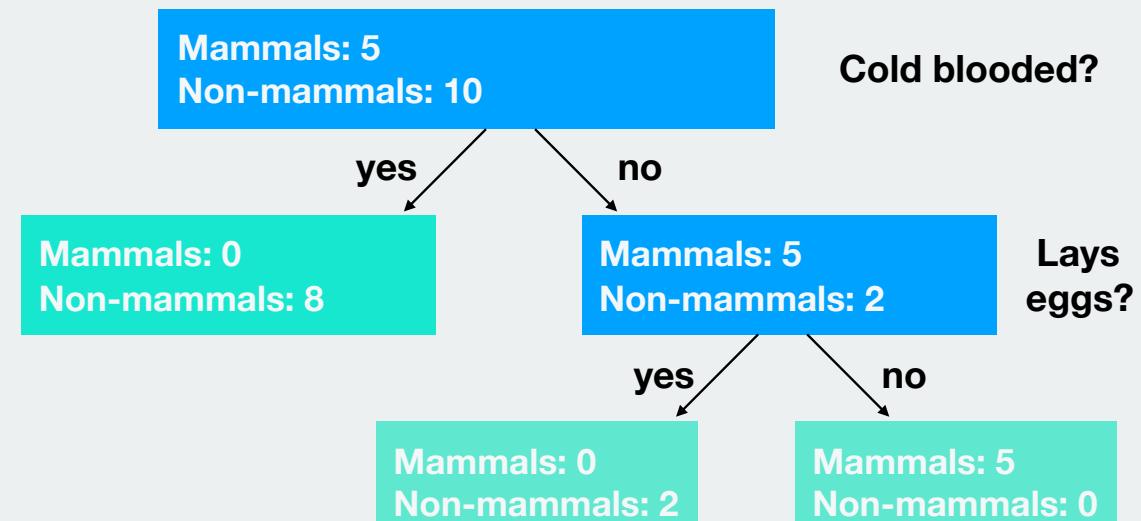
Decision trees

Could we have asked better questions?

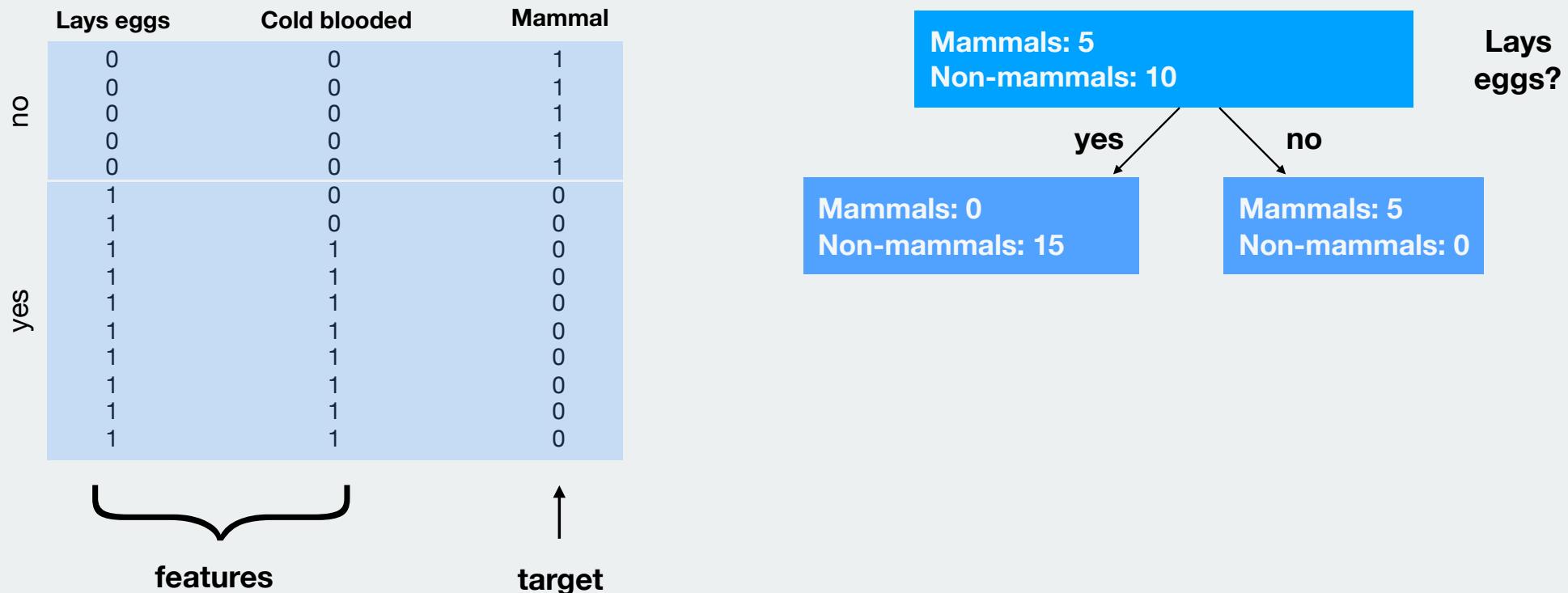
Decision trees

	Lays eggs	Cold blooded	Mammal
no	0 0 0 0 0 0	0 0 0 0 0 0	1 1 1 1 1 1
yes	1 1	0 0	0 0
	1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1	0 0 0 0 0 0 0 0 0 0

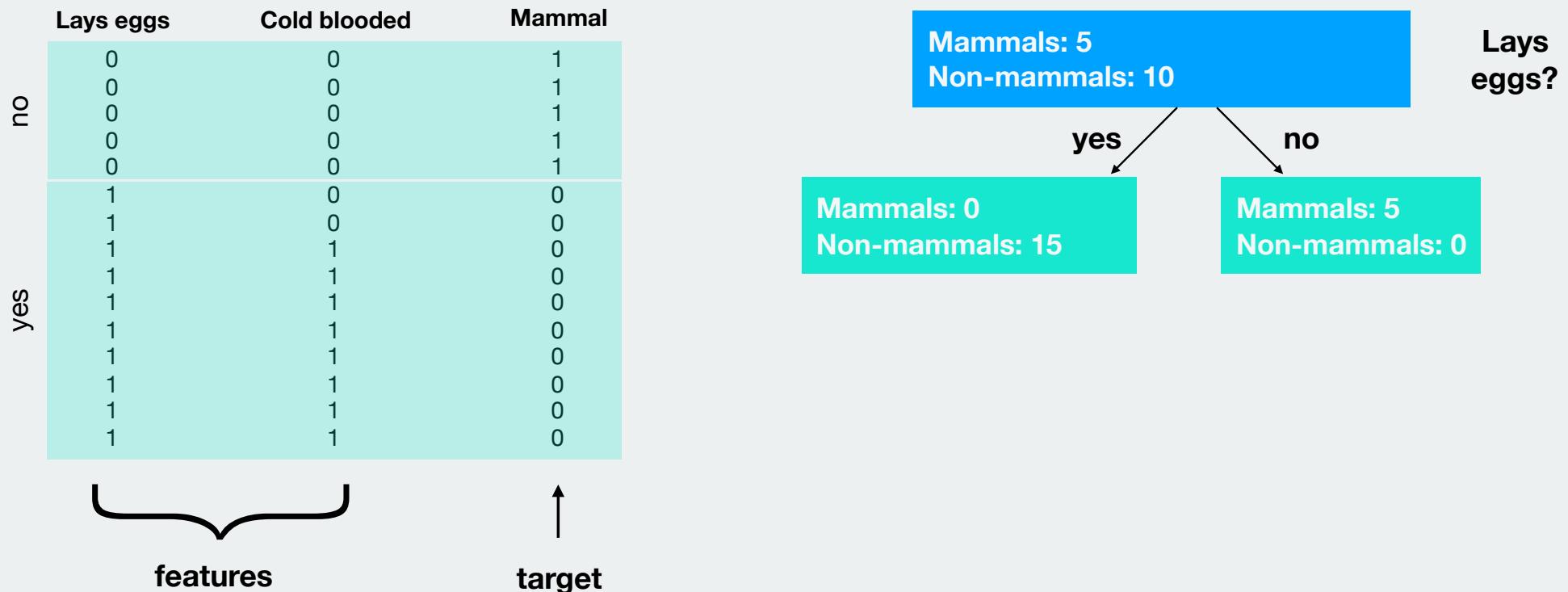
features **target**



Decision trees

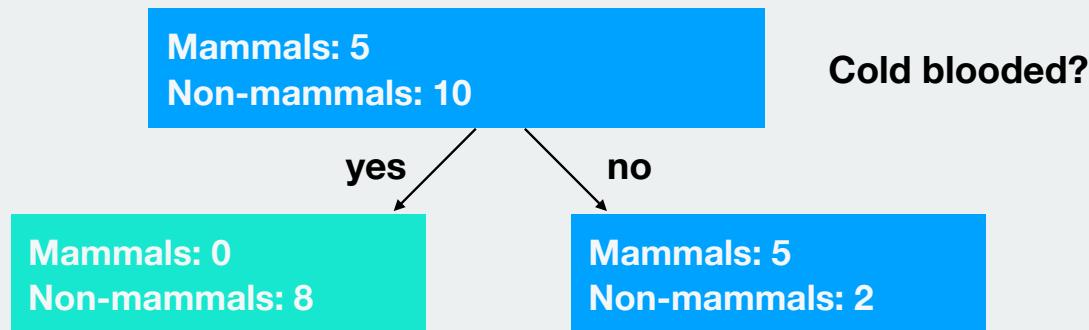


Decision trees

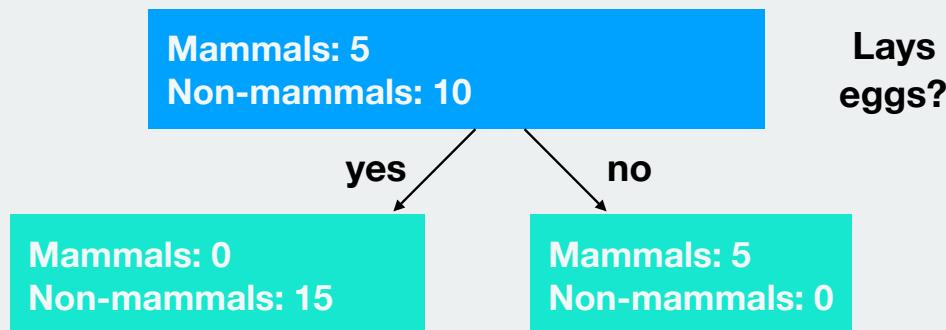


Split selection

Split 1:



Split 2:



Split selection

$$\text{(Shannon) } Entropy = - \sum_i p(i) \log_2 p(i)$$

Input: Probability vector (a list of values between 0 and 1, which sums to 1)

Output: Entropy (a measure of how “spread out” the probability distribution is)

Split selection

$$\text{Entropy} = - \sum_i p(i) \log_2 p(i)$$

Mammals: 0
Non-mammals: 8

$$p = [1, 0]$$

Split selection

$$\text{Entropy} = - \sum_i p(i) \log_2 p(i)$$

Mammals: 0
Non-mammals: 8

$$p = [1, 0]$$

$$\text{Entropy} = - (1 \cdot \log_2(1) + 0 \cdot \log_2(0)) = 0$$

Split selection

$$\text{Entropy} = - \sum_i p(i) \log_2 p(i)$$

Mammals: 0
Non-mammals: 8

$$p = [1, 0]$$

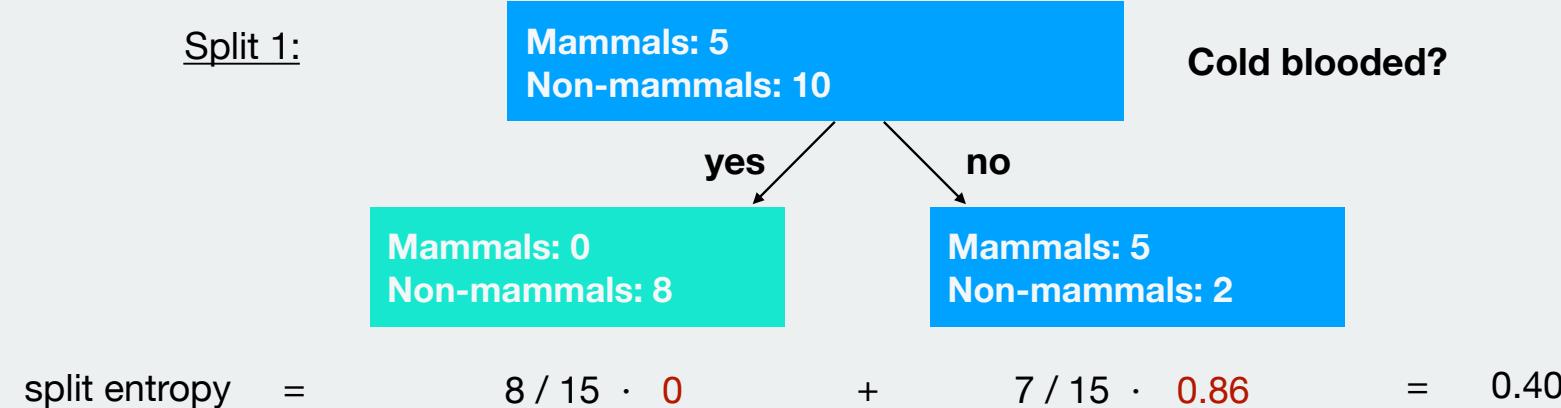
$$\text{Entropy} = - (1 \cdot \log_2(1) + 0 \cdot \log_2(0)) = 0$$

Mammals: 5
Non-mammals: 2

$$p = [2/7, 5/7]$$

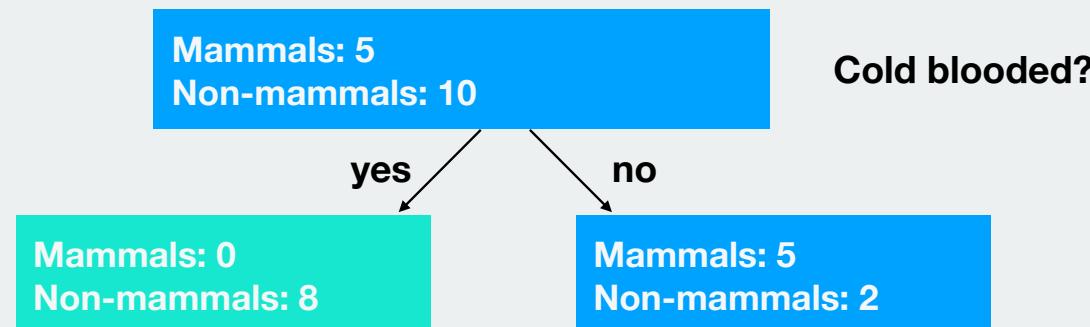
$$\text{Entropy} = - (2/7 \cdot \log_2(2/7) + 5/7 \cdot \log_2(5/7)) = 0.86$$

Split selection



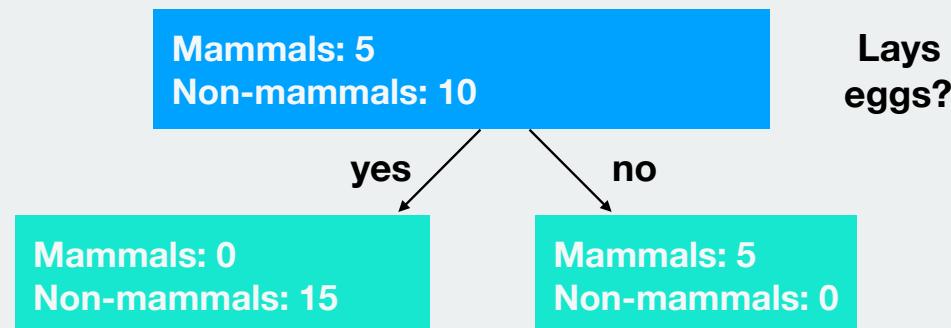
Automatic split selection

Split 1:



$$\text{split entropy} = \frac{8}{15} \cdot 0 + \frac{7}{15} \cdot 0.86 = 0.40$$

Split 2:

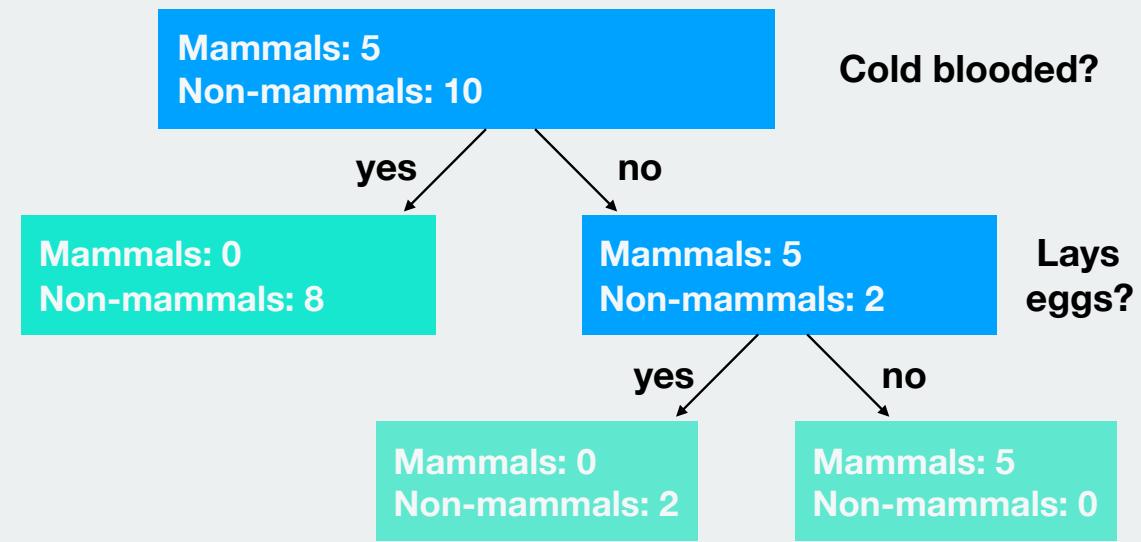


$$\text{split entropy} = \frac{8}{15} \cdot 0 + \frac{7}{15} \cdot 0 = 0$$

Decision trees

	Lays eggs	Cold blooded	Mammal
no	0 0 0 0 0 0	0 0 0 0 0 0	1 1 1 1 1 1
yes	1 1	0 0	0 0
	1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1	0 0 0 0 0 0 0 0 0 0

features
target



Decision trees

- Each node is split based on the attribute with the **best** split entropy
- Leafs = decisions (majority vote)
- To classify a new input, start at the root and make your way down through the tree by answering attribute questions at each node

Decision trees

What do we do when we have Big(ger) Data?

Decision trees

	A-Force	A-Next	A.R.M.O.R.	Acolytes	Activation Pack	Advanced Idea.	Age of Apocal.	Agency X	...	Young Avengers	Young X-Men	Zodiac
Abigail Brand	0	0	1	0	0	0	0	1	...	1	0	0
Adam Warlock	0	1	0	0	0	0	0	0	...	0	1	0
Adept (comics)	0	0	0	0	0	1	0	0	...	0	0	0
...									...			
Zzzax	1	0	0	0	0	0	1	0	...	0	0	0

Decision trees ● ●

Ensemble learning ○

Logistic regression ○

On steroids:
Ensemble Learning

Ensemble Learning

- Create and train many classification models
- Treat each model as a “voter”
- For each datapoint, classify it according to what models predicts it to be

Random Forest

model1

	Pclass1	Pclass2	Pclass3	Sexfemale	Sexmale	Embarkedn	EmbarkedC	EmbarkedQ	EmbarkedS	CabinFalse	CabinTrue	PassengerId	Age	SibSp	Parch	Fare	Survived
0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	1	22.0	1	0	7.2500	0
1	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	2	38.0	1	0	71.2833	1
2	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	3	26.0	0	0	7.9250	1
3	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	4	35.0	1	0	53.1000	1
4	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	5	35.0	0	0	8.0500	0
5	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	0.0	6	Nan	0	0	8.4583	0
6	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	7	54.0	0	0	51.8625	0
7	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	8	2.0	3	1	21.0750	0
8	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	9	27.0	0	2	11.1333	1
9	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	10	14.0	1	0	30.0708	1
...
881	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	882	33.0	0	0	7.8958	0
882	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	883	22.0	0	0	10.5167	0
883	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	884	28.0	0	0	10.5000	0
884	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	885	25.0	0	0	7.0500	0
885	0.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	886	39.0	0	5	29.1250	0
886	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	887	27.0	0	0	13.0000	0
887	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	888	19.0	0	0	30.0000	1
888	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	889	Nan	1	2	23.4500	0
889	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	890	26.0	0	0	30.0000	1
890	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	0.0	891	32.0	0	0	7.7500	0

Random Forest

model2

	Pclass1	Pclass2	Pclass3	Sexfemale	Sexmale	Embarkednan	EmbarkedC	EmbarkedQ	EmbarkedS	CabinFalse	CabinTrue	PassengerId	Age	SibSp	Parch	Fare	Survived
0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	1	22.0	1	0	7.2500	0
1	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	2	38.0	1	0	71.2833	1
2	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	3	26.0	0	0	7.9250	1
3	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	4	35.0	1	0	53.1000	1
4	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	5	35.0	0	0	8.0500	0
5	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	0.0	6	NaN	0	0	8.4583	0
6	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	7	54.0	0	0	51.8625	0
7	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	8	2.0	3	1	21.0750	0
8	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	9	27.0	0	2	11.1333	1
9	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	10	14.0	1	0	30.0708	1
...
881	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	882	33.0	0	0	7.8958	0
882	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	883	22.0	0	0	10.5167	0
883	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	884	28.0	0	0	10.5000	0
884	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	885	25.0	0	0	7.0500	0
885	0.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	886	39.0	0	5	29.1250	0
886	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	887	27.0	0	0	13.0000	0
887	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	888	19.0	0	0	30.0000	1
888	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0	889	NaN	1	2	23.4500	0
889	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	890	26.0	0	0	30.0000	1
890	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	0.0	891	32.0	0	0	7.7500	0

Random Forest

model3

Pclass1	Pclass2	Pclass3	Sexfemale	Sexmale	Embarkednan	EmbarkedC	EmbarkedQ	EmbarkedS	CabinFalse	CabinTrue	PassengerId	Age	SibSp	Parch	Fare	Survived
0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	1	22.0	1	0	7.2500	0
1	1.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0	2	38.0	1	0	71.2833	1
2	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	3	26.0	0	0	7.9250	1
3	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	4	35.0	1	0	53.1000	1
4	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	5	35.0	0	0	8.0500	0
5	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	6	Nan	0	0	8.4583	0
6	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	7	54.0	0	0	51.8625	0
7	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	8	2.0	3	1	21.0750	0
8	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0	9	27.0	0	2	11.1333	1
9	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0	10	14.0	1	0	30.0708	1
...
881	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	1.0	0.0	882	33.0	0	0	7.8958	0
882	0.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0	1.0	0.0	883	22.0	0	0	10.5167	0
883	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	884	28.0	0	0	10.5000	0
884	0.0	0.0	1.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	885	25.0	0	0	7.0500	0
885	0.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	886	39.0	0	5	29.1250	0
886	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	1.0	887	27.0	0	0	13.0000	0
887	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	888	19.0	0	0	30.0000	1
888	0.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	889	Nan	1	2	23.4500	0
889	1.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	890	26.0	0	0	30.0000	1
890	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0	0.0	1.0	891	32.0	0	0	7.7500	0

Ensemble Learning

- Create and train many classification models
- Treat each model as a “voter”
- For each datapoint, classify it according to what models predicts it to be

$\text{model1}(x) = 1$
 $\text{model2}(x) = 1$
 $\text{model3}(x) = 0$
 $\text{model4}(x) = 1$
 $\text{model5}(x) = 1$
 $\text{model6}(x) = 1$
 $\text{model7}(x) = 0$
 $\text{model8}(x) = 1$
...
 $\text{modeln}(x) = 1$

average = 0.84 ≈ 1

Ensemble Learning – benefits

- Very fast training (fewer attributes)
- Pretty good performance
- Very robust to overfitting

So far today

- Decision Trees (DTs) = a ML algorithm
 - Entropy, split decisions
- Ensemble Learning (EL) = using many models (+/- many algorithms)
 - Random Forrest, an example of EL using DTs

Decision trees ● ●

Ensemble learning ●

Logistic regression ○

Linear (and Logistic) regression

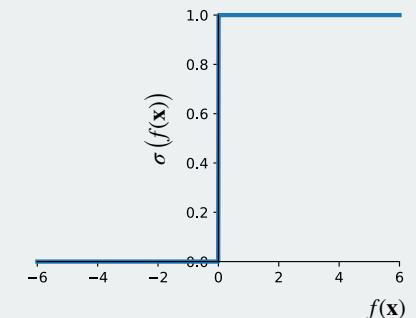
Linear regression

$$w_0 + x_0 w_1 + x_1 w_2 + x_2 w_3 = f(\mathbf{x})$$

Linear regression classifier

$$w_0 + x_0 w_1 + x_1 w_2 + x_2 w_3 = f(\mathbf{x})$$

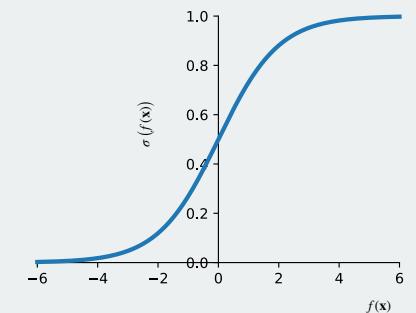
$$\sigma(f(\mathbf{x})) = \begin{cases} 1, & \text{if } f(\mathbf{x}) \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



Linear regression regression

$$w_0 + x_0 w_1 + x_1 w_2 + x_2 w_3 = f(\mathbf{x})$$

$$\sigma(f(\mathbf{x})) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

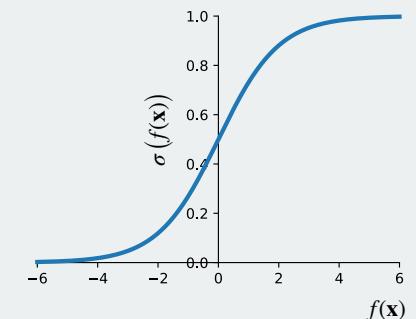


Linear regression classifier

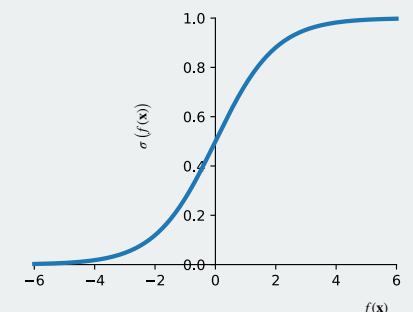
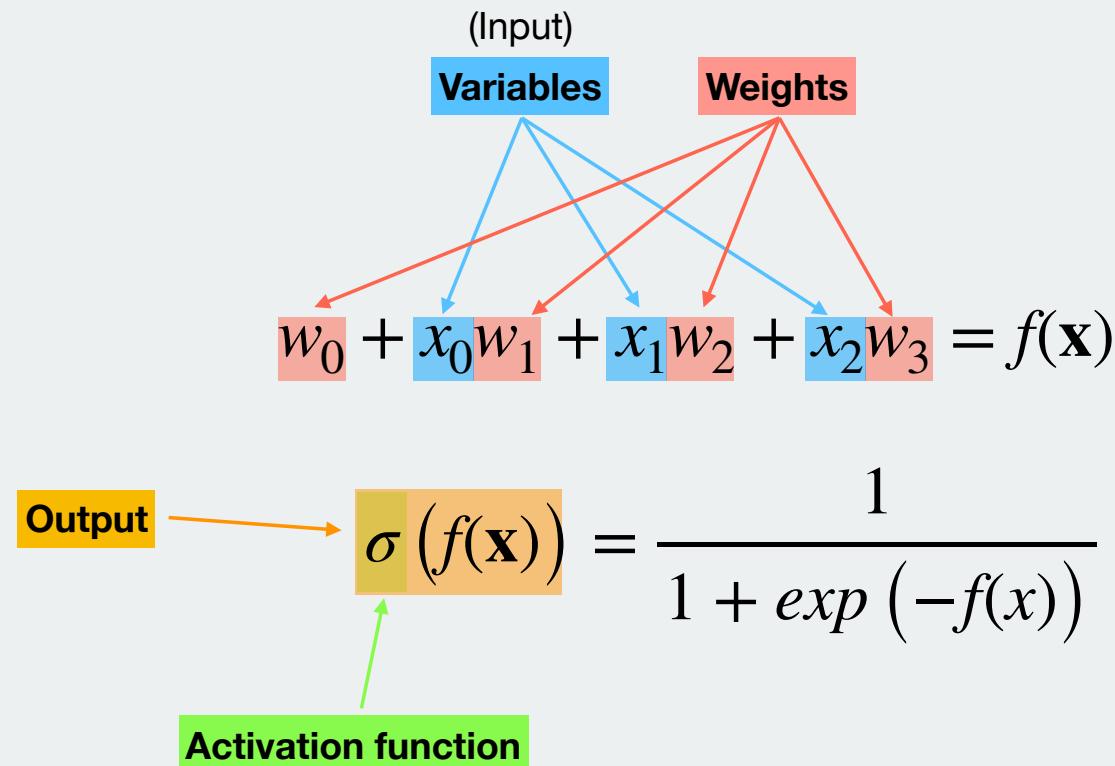
*i*Logistic regression!
classifier

$$w_0 + x_0 w_1 + x_1 w_2 + x_2 w_3 = f(\mathbf{x})$$

$$\sigma(f(\mathbf{x})) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$



Logistic regression classifier



Sum up

- Decision Trees (DTs) = a ML algorithm
 - Entropy, split decisions
- Ensemble Learning (EL) = using many models (+/- many algorithms)
 - Random Forrest, an example of EL using DTs
- Linear/Logistic Regression = a ML algorithm
 - Linear regression outputs 0 or 1
 - Logistic regression gives a probability