

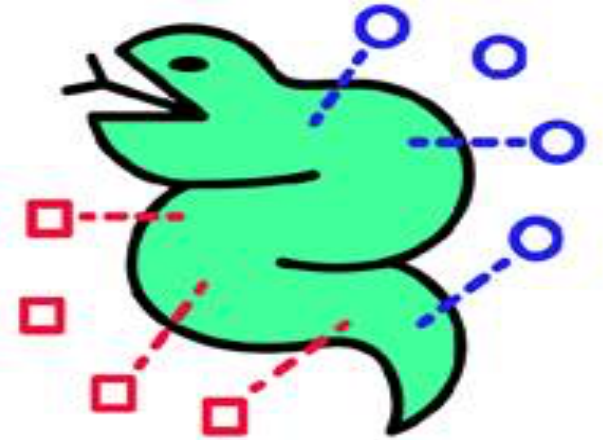
Decision Trees

Alexander Jung

Assistant Professor

Department of Computer Science

Aalto University



Machine Learning
With Python

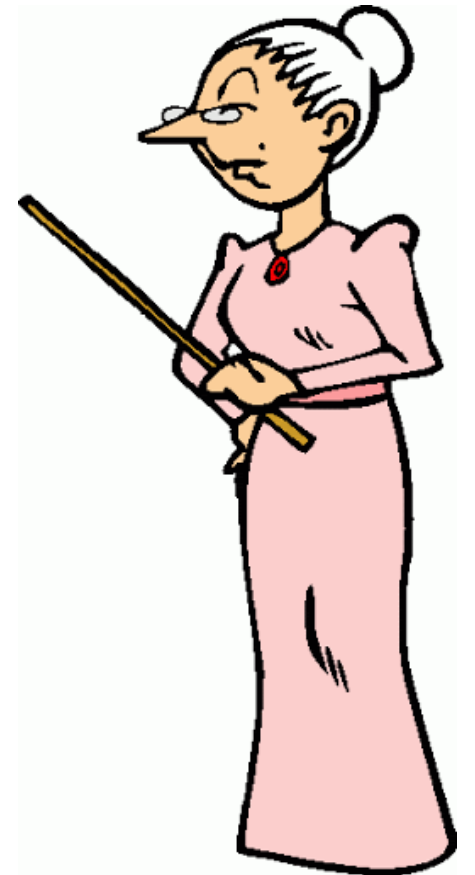


FITECH
NETWORK UNIVERSITY



Aalto-yliopisto

What are three **main**
components of machine
learning?



1. Data

Data

- set of “data points” (atomic unit of information)
- data point has **features and labels**
- **features** are properties that **can measured easily**
- **labels** =higher-level facts or quantities of interest

Data Point = "Some Movie"

features:

- x_1 = movie duration in minutes
- x_2 = screen time of **Arnold Schwarzenegger**

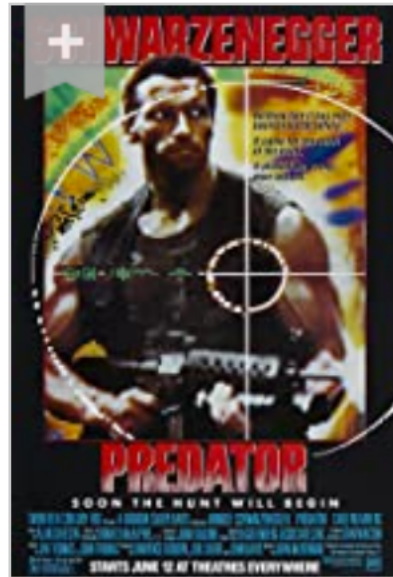


label: the movie rating (scale 0 ... 10)

Data = Bunch of Data Points



The 6th Day
Adam Gibson
(2000)



Predator
Dutch
(1987)



Eraser
U.S. Marshal John 'The E...'
(1996)



Last Action Hero
Jack Slater
(1993)



Scatter Plot

x2



○ "Terminator 1"

○ "Twins"

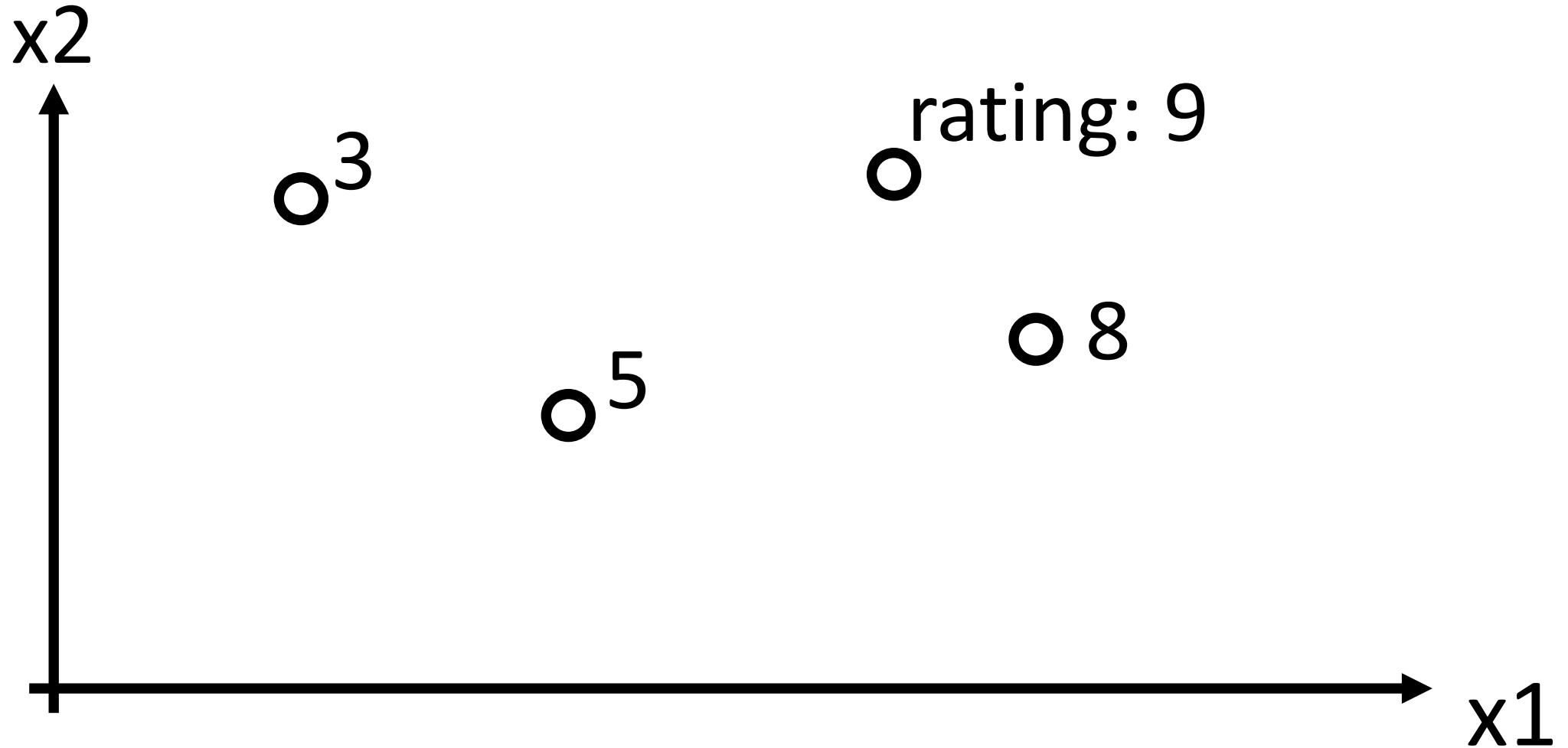


"The Expendables 3"



x1

Scatter Plot

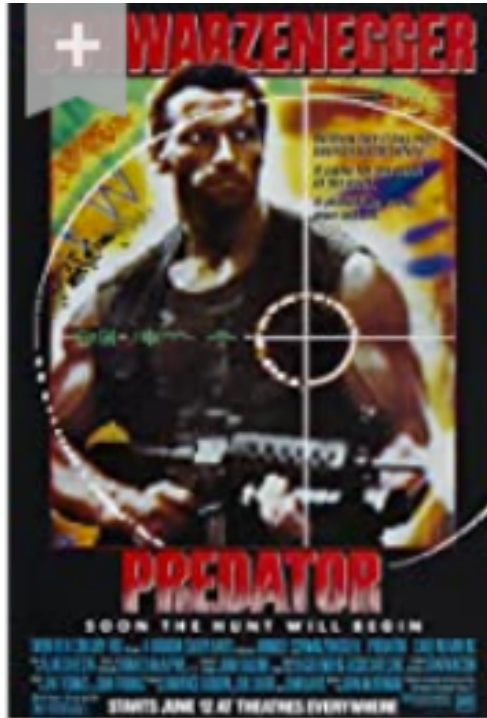


Different Choice of Label

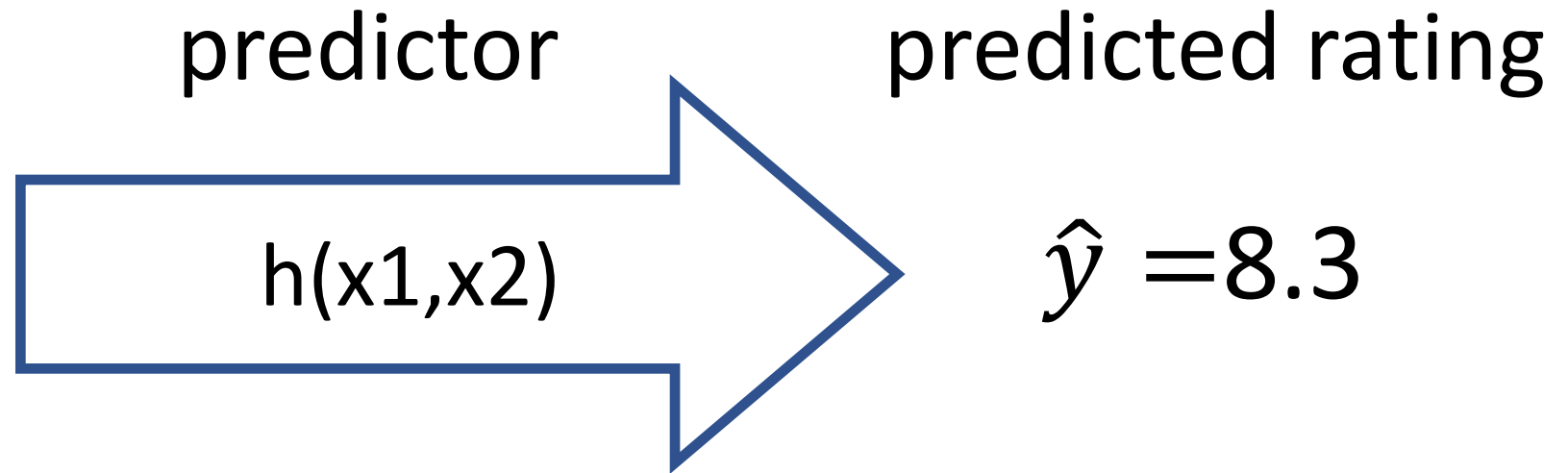


2. Hypothesis Space

How Many Predictors Are There?



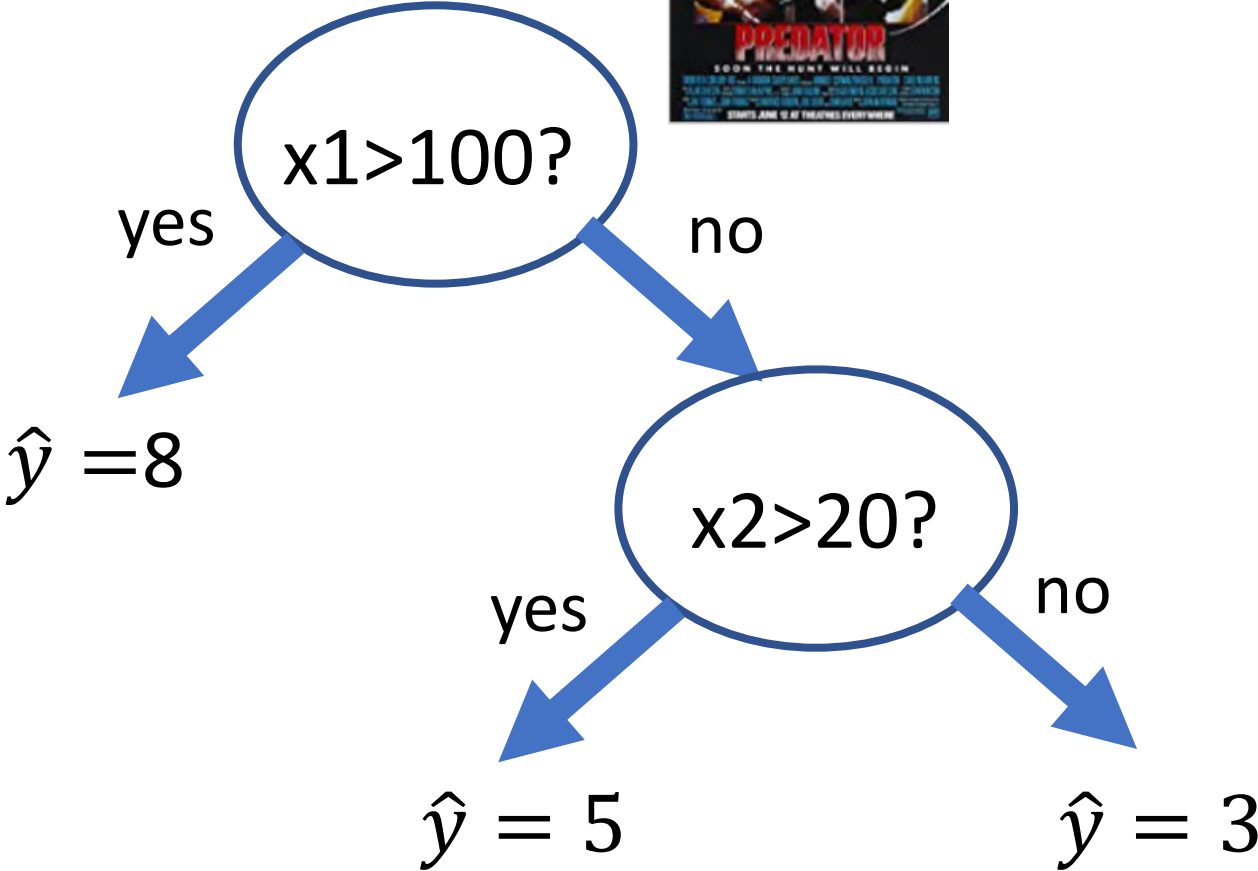
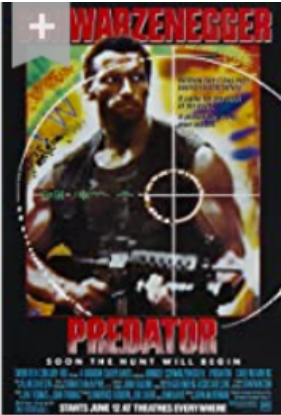
$x_1 = 123.348$,
 $x_2 = 40.456$



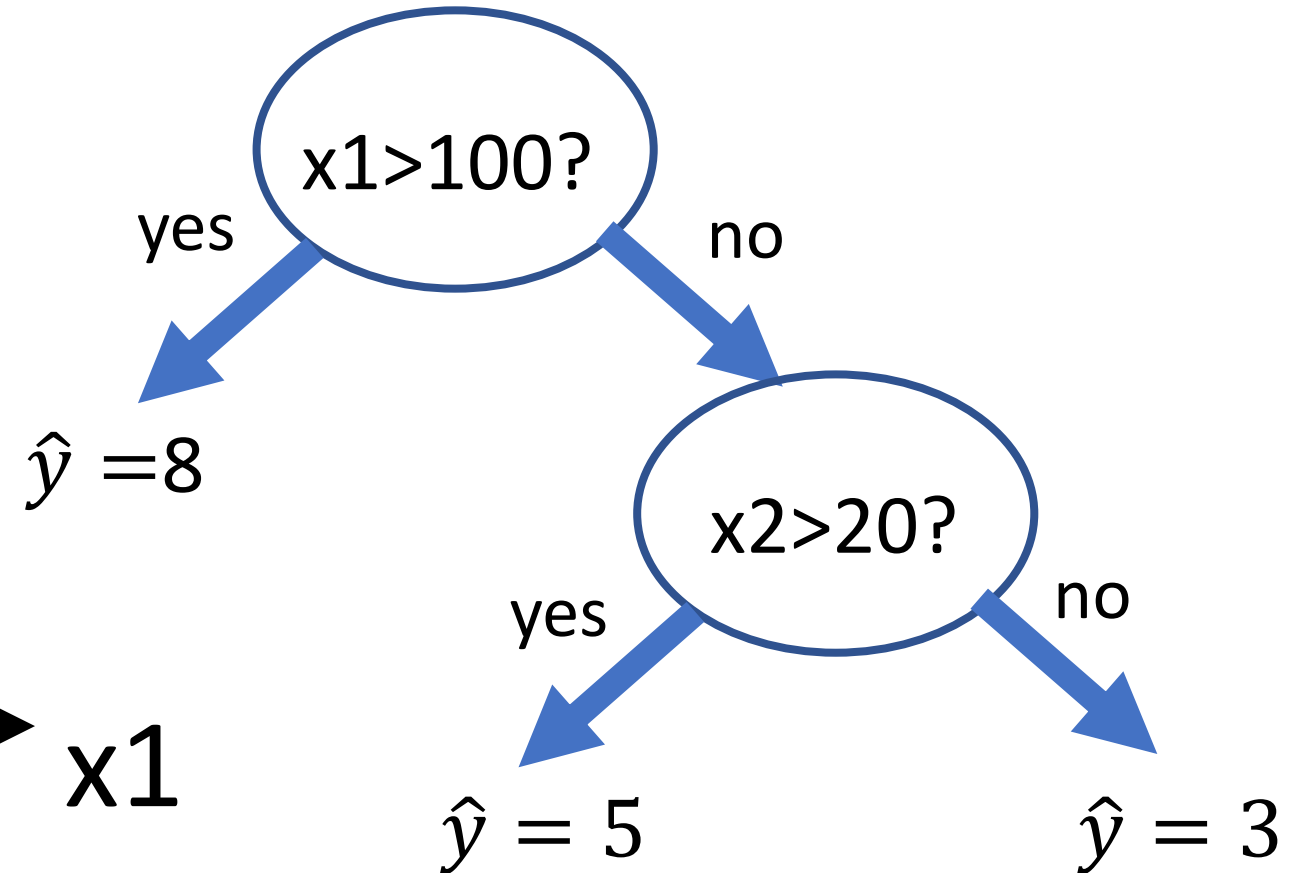
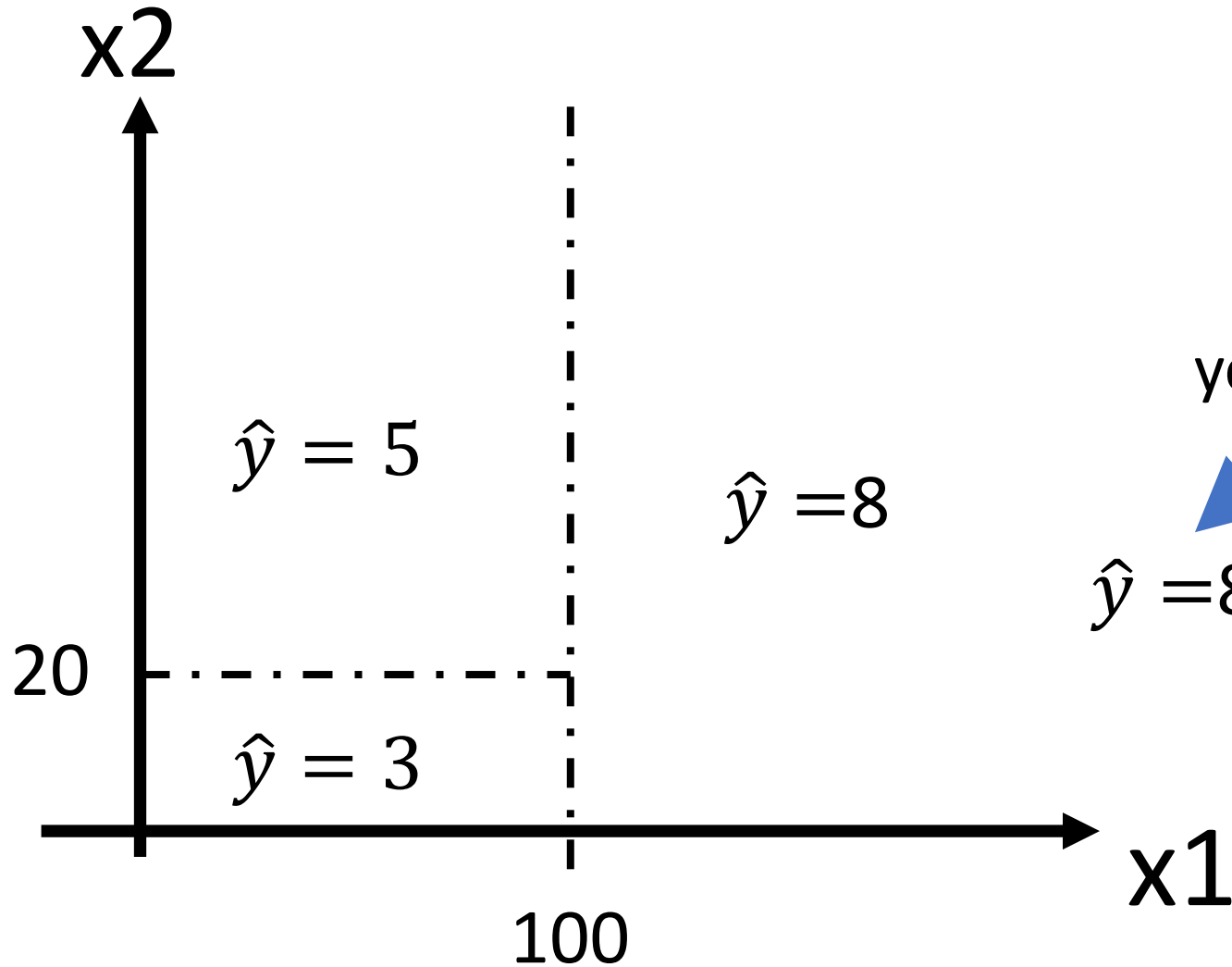
Hypothesis Space “Decision Tree”

- infinitely many functions from x_1, x_2 to \hat{y}
- restrict to subset of maps (hypothesis space)
- subset of maps given by decision trees

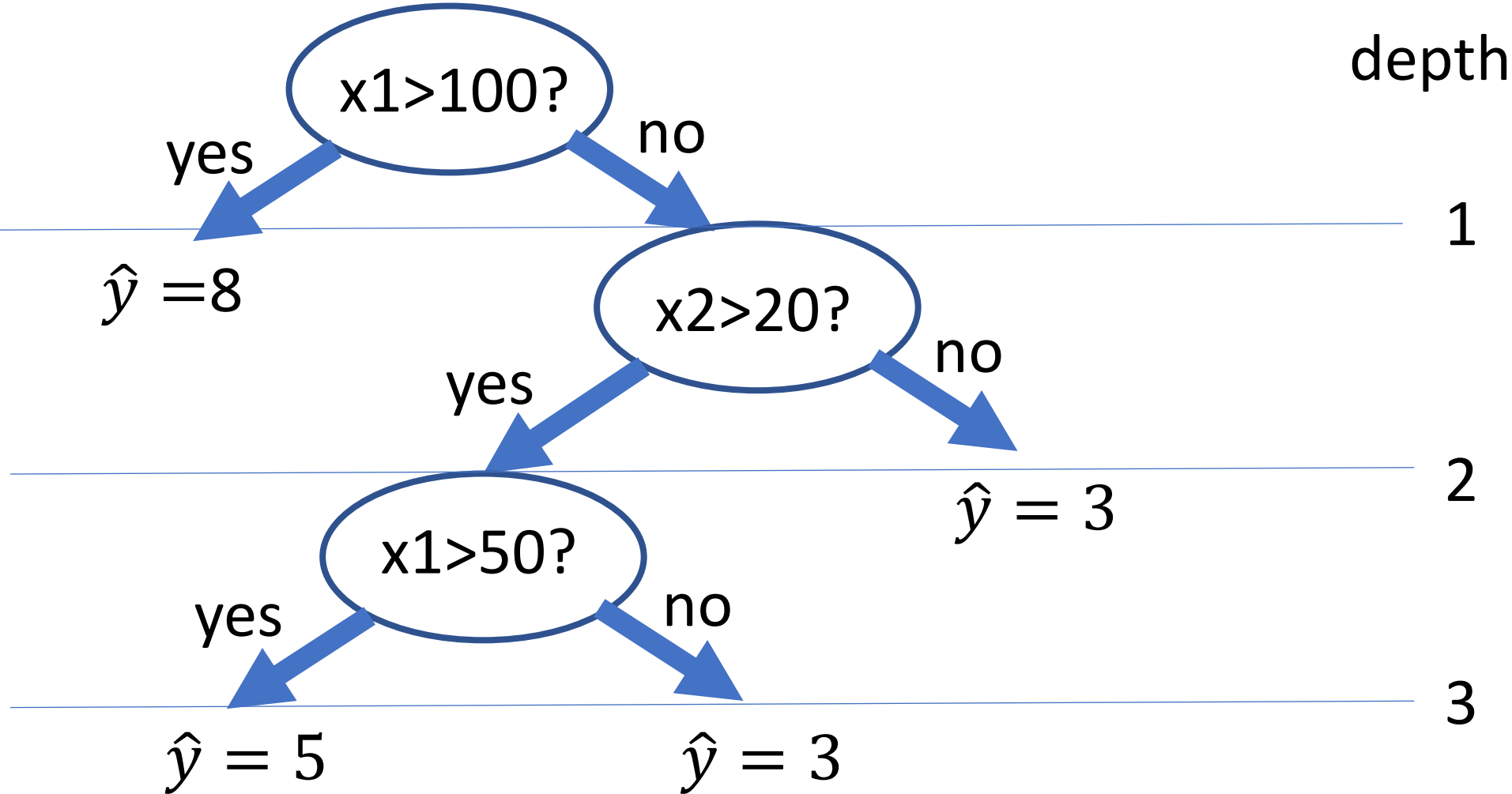
A Decision Tree (DT)



DT in Feature Space



Tree Depth



Hypothesis Space of DTs

space of predictor maps given by DT that involve threshold tests for x_1 and x_2

`sklearn.tree.DecisionTreeRegressor`

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='mse', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, presort='deprecated', ccp_alpha=0.0) \[source\]
```

A decision tree regressor.

numeric labels

`sklearn.tree.DecisionTreeClassifier`

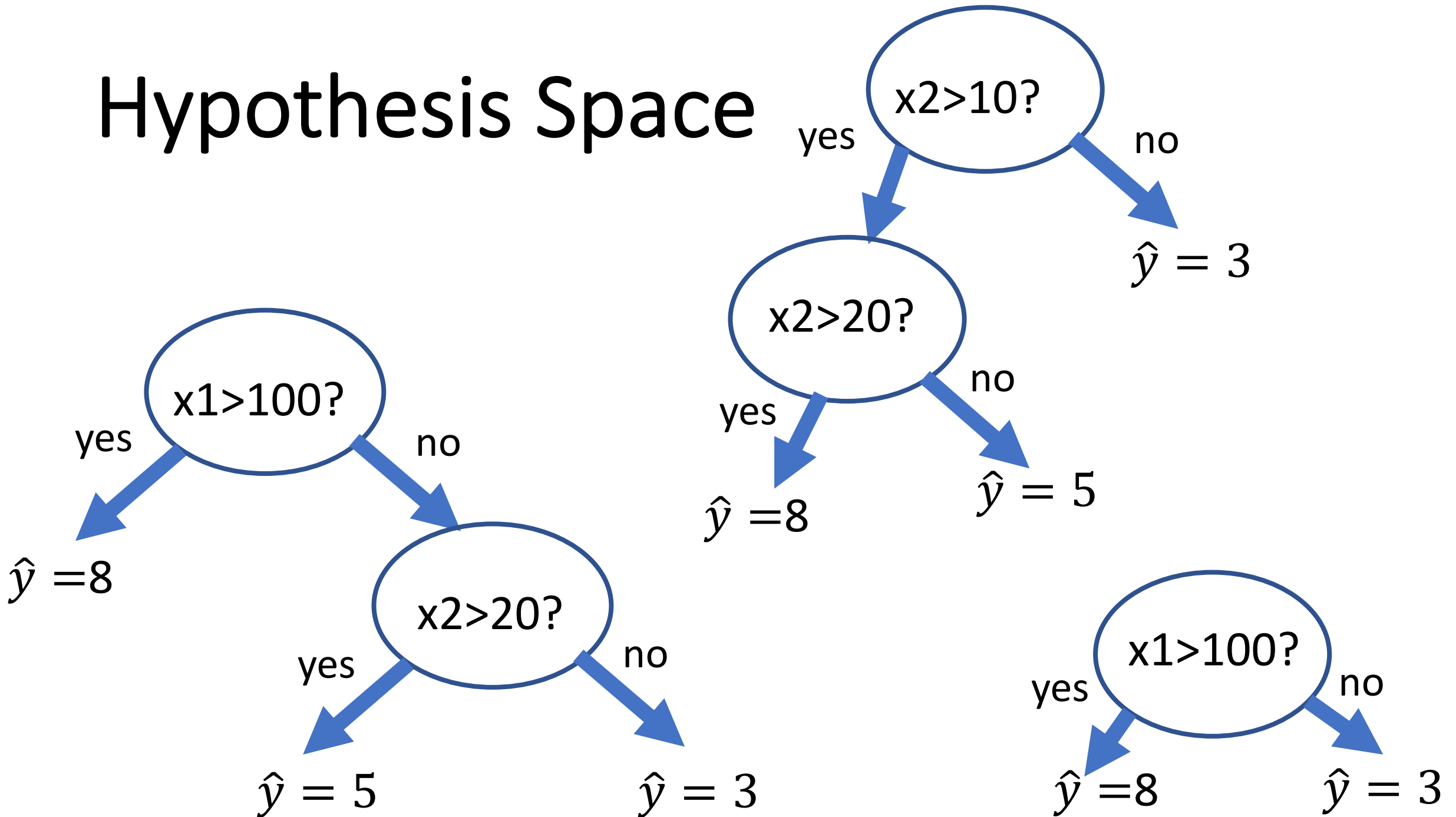
```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, class_weight=None, presort='deprecated', ccp_alpha=0.0) \[source\]
```

A decision tree classifier.

Read more in the [User Guide](#).

discrete valued
label

Hypothesis Space



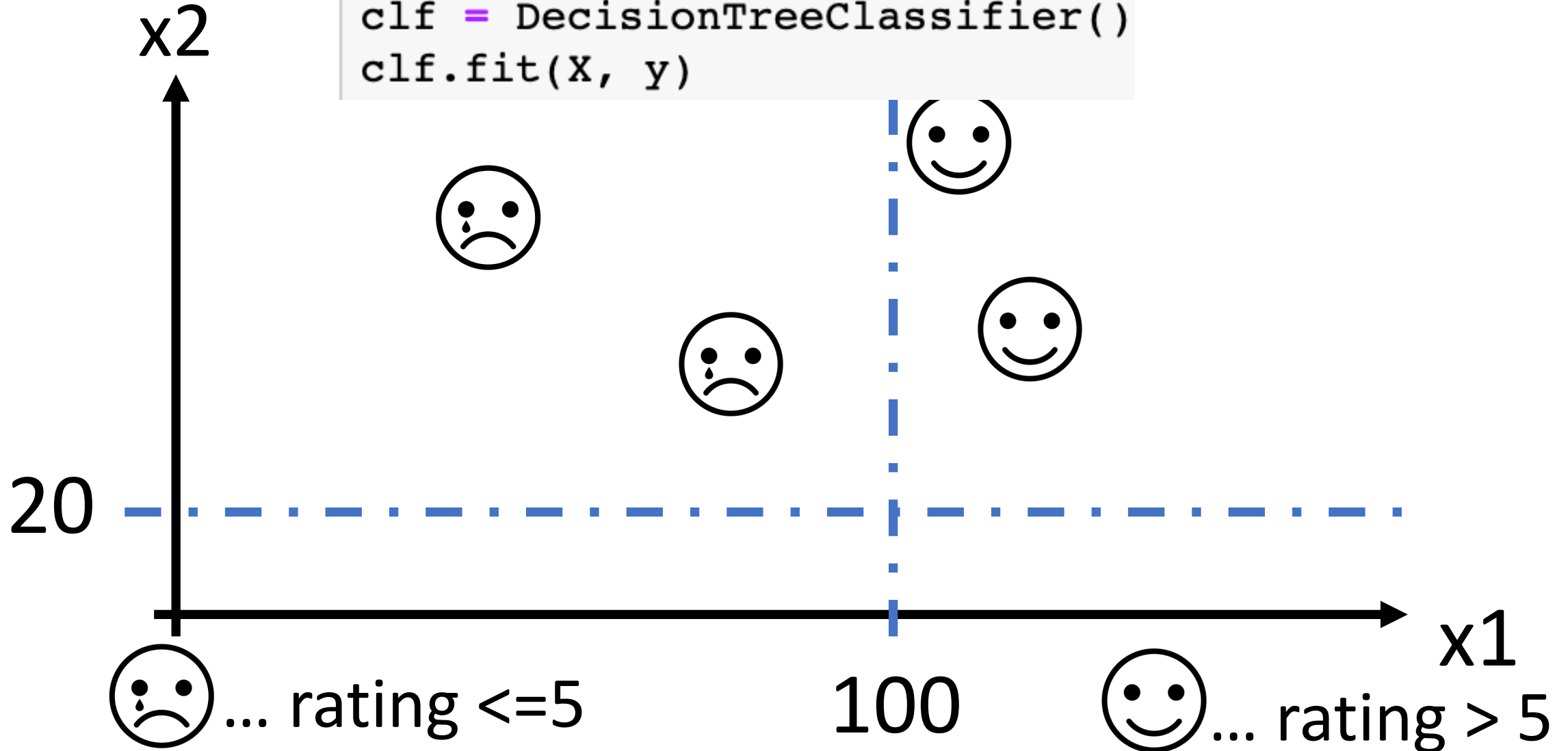
3. Loss Function

Learning a Good DT

- decision tree defines a hypothesis space
- set of maps that are represented by DT
- quality of map measured by (average) loss
- can use **any loss function**

Some Labeled (Training) Data

```
clf = DecisionTreeClassifier()  
clf.fit(X, y)
```

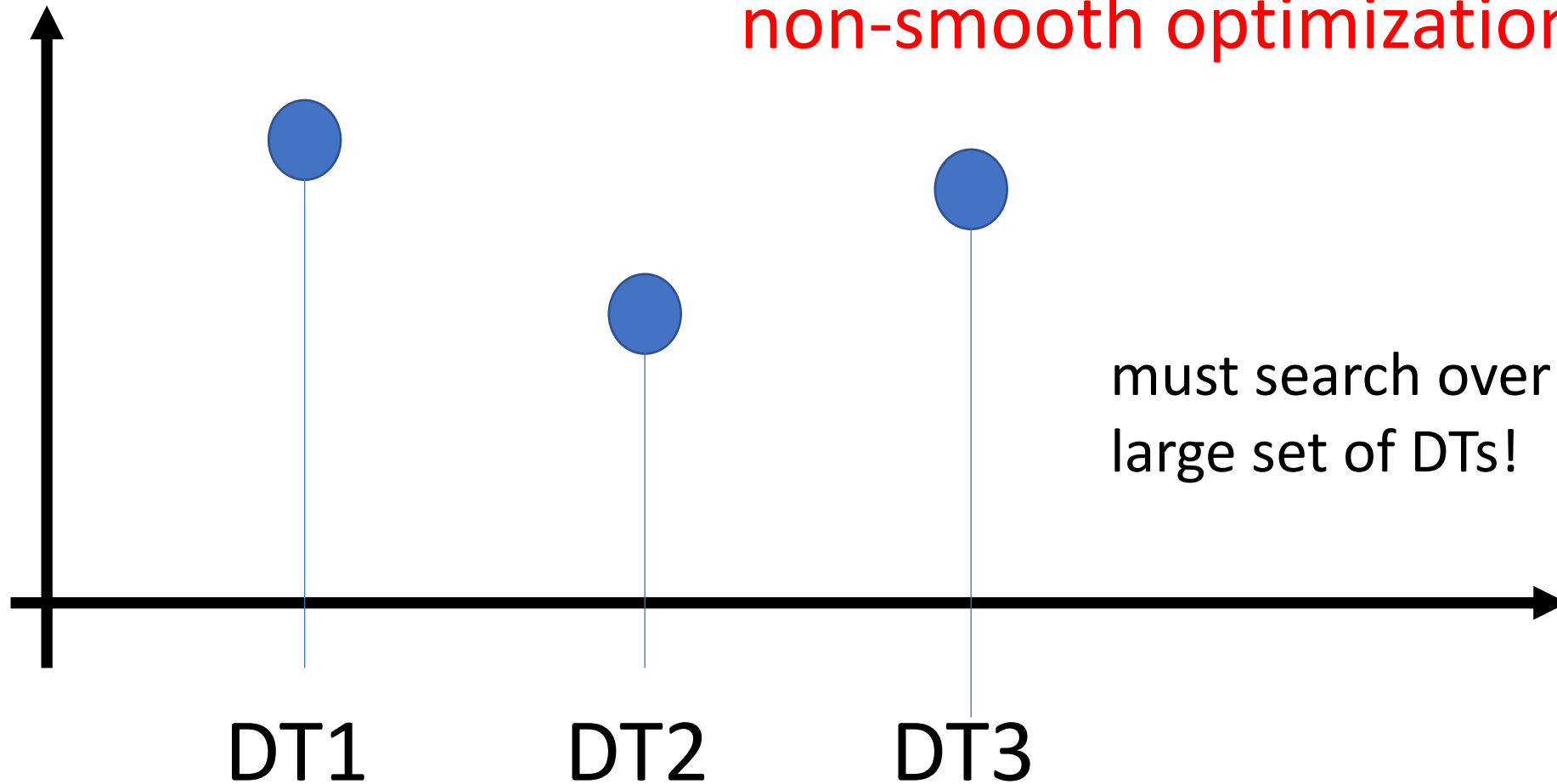


Loss Minimization



non-convex,
non-smooth optimization problem

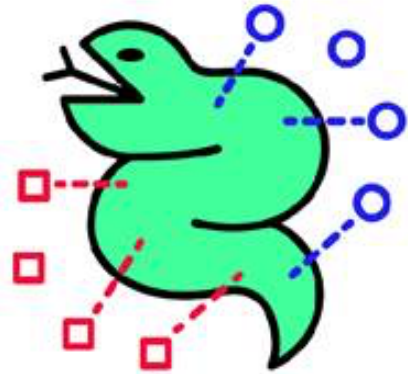
average loss



must search over
large set of DTs!

So What ?

- DT is a flow chart of predictor map
- DT define a hypothesis space
- DT can be combined with different loss functions
- DT can be used for regression or classification



Machine Learning
With Python



Thank You !