# Decision Trees

# Announcements

- No class on Wednesday Jan 29th - I'm out of town.

- Assignment still due on that Wed afternoon. - you have all information you need.

- Office hours: Wednesdays, 3-5PM (I'll be there today, but not next week.)

# Decision Trees

---

**Algorithm 1** DECISIONTREETRAIN(*data*, *remaining features*)

---

1:  *guess* ← most frequent answer in *data*          // default answer for this data
2:  **if** the labels in *data* are unambiguous **then**
3:     **return** LEAF(*guess*)               // base case: no need to split further
4:  **else if** *remaining features* is empty **then**
5:     **return** LEAF(*guess*)               // base case: cannot split further
6:  **else**                                   // we need to query more features
7:     **for all** $f \in$ *remaining features* **do**
8:        $NO$ ← the subset of *data* on which $f$=*no*
9:        $YES$ ← the subset of *data* on which $f$=*yes*
10:      $score[f]$ ← # of majority vote answers in $NO$
11:                 + # of majority vote answers in $YES$
                             // the accuracy we would get if we only queried on *f*
12:     **end for**
13:     $f$ ← the feature with maximal $score(f)$
14:     $NO$ ← the subset of *data* on which $f$=*no*
15:     $YES$ ← the subset of *data* on which $f$=*yes*
16:     *left* ← DECISIONTREETRAIN(*NO*, *remaining features* \ {$f$})
17:     *right* ← DECISIONTREETRAIN(*YES*, *remaining features* \ {$f$})
18:     **return** NODE(*f*, *left*, *right*)
19: **end if**

---

# Decision Trees

---

**Algorithm 2** DECISIONTREETEST(*tree*, *test point*)

---

1: **if** *tree* is of the form LEAF(*guess*) **then**
2:     **return** *guess*
3: **else if** *tree* is of the form NODE(*f*, *left*, *right*) **then**
4:     **if** $f = no$ in *test point* **then**
5:         **return** DECISIONTREETEST(*left*, *test point*)
6:     **else**
7:         **return** DECISIONTREETEST(*right*, *test point*)
8:     **end if**
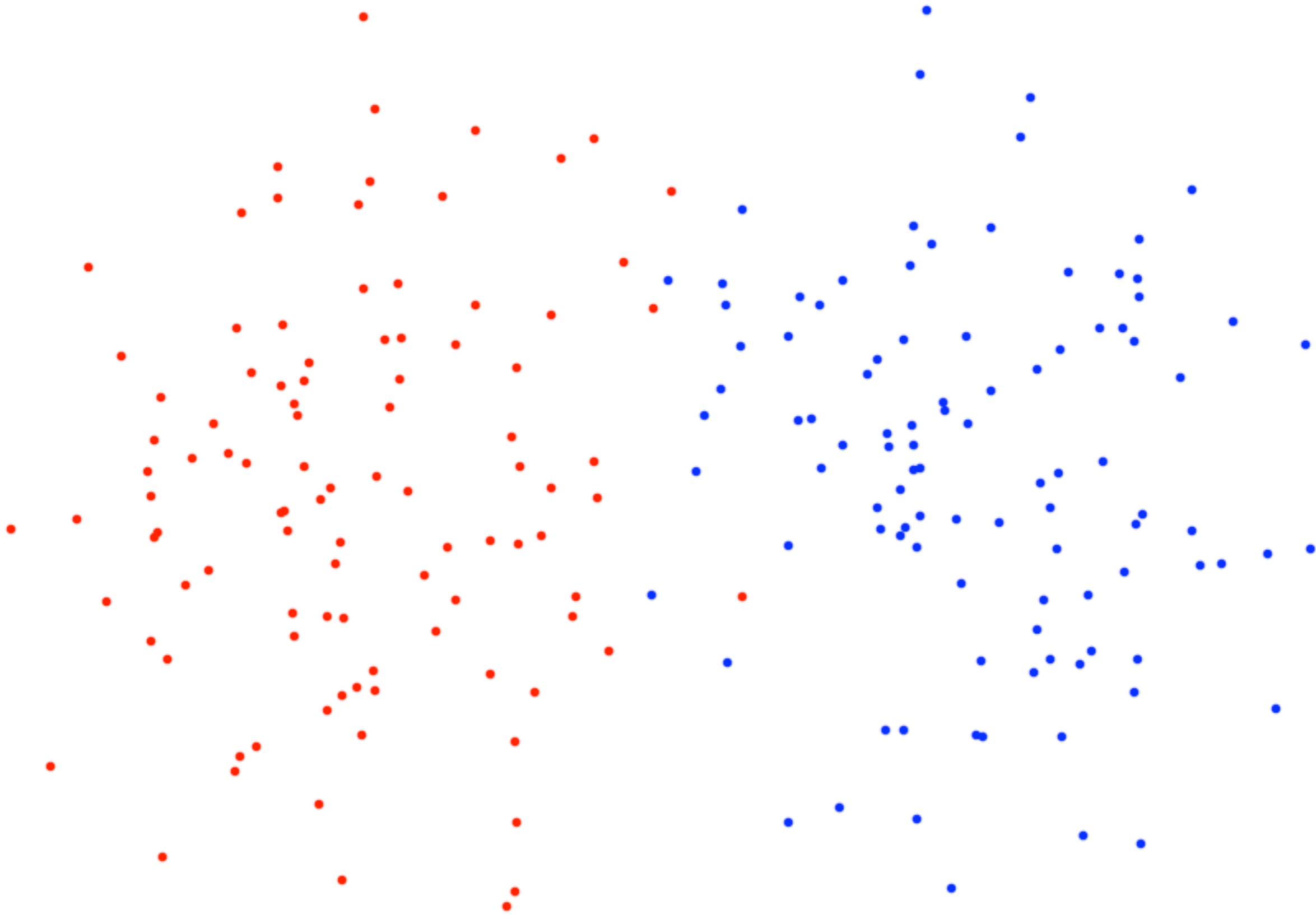9: **end if**

---

# Questions
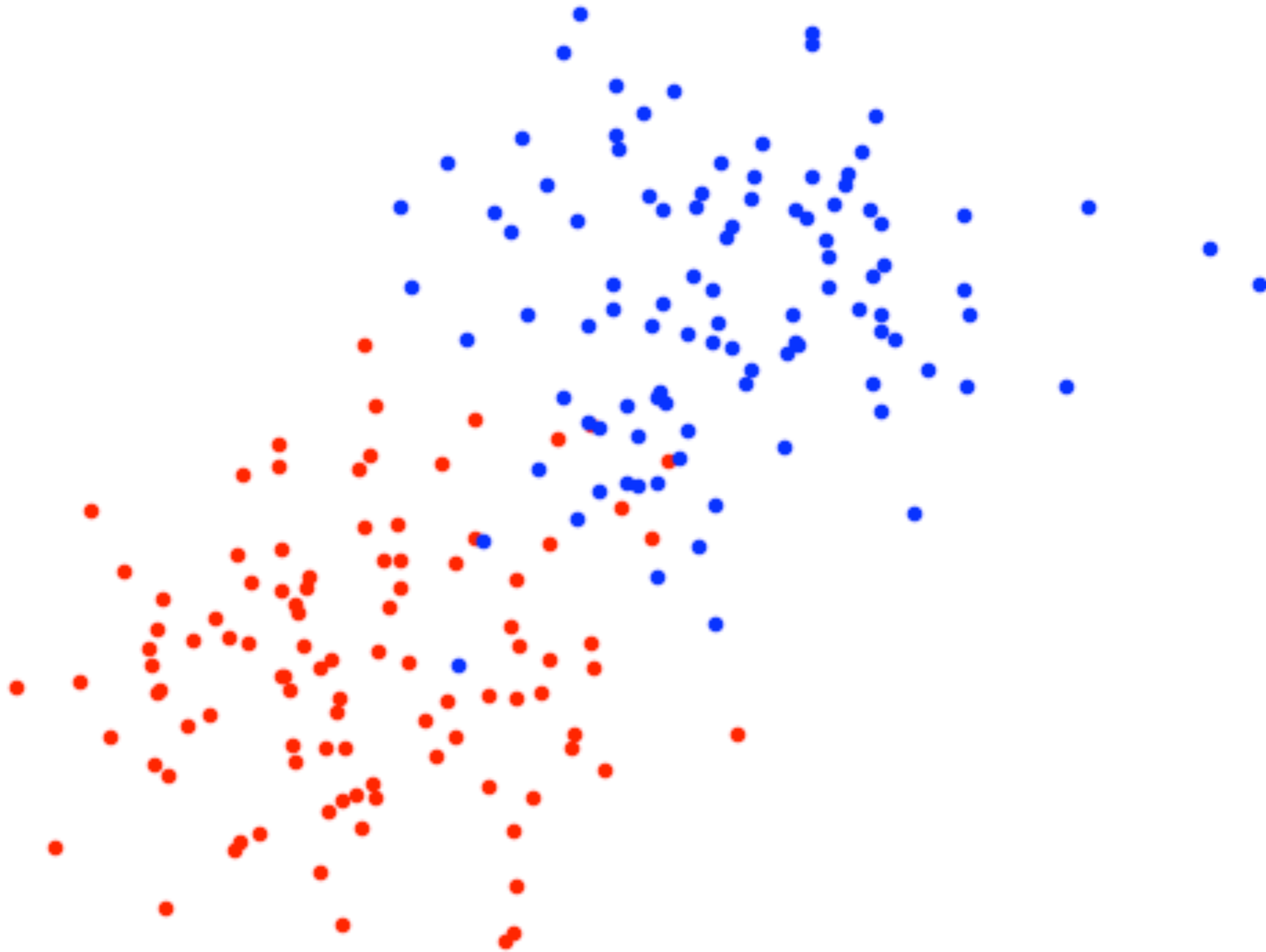
- Is this greedy strategy good?

# Questions

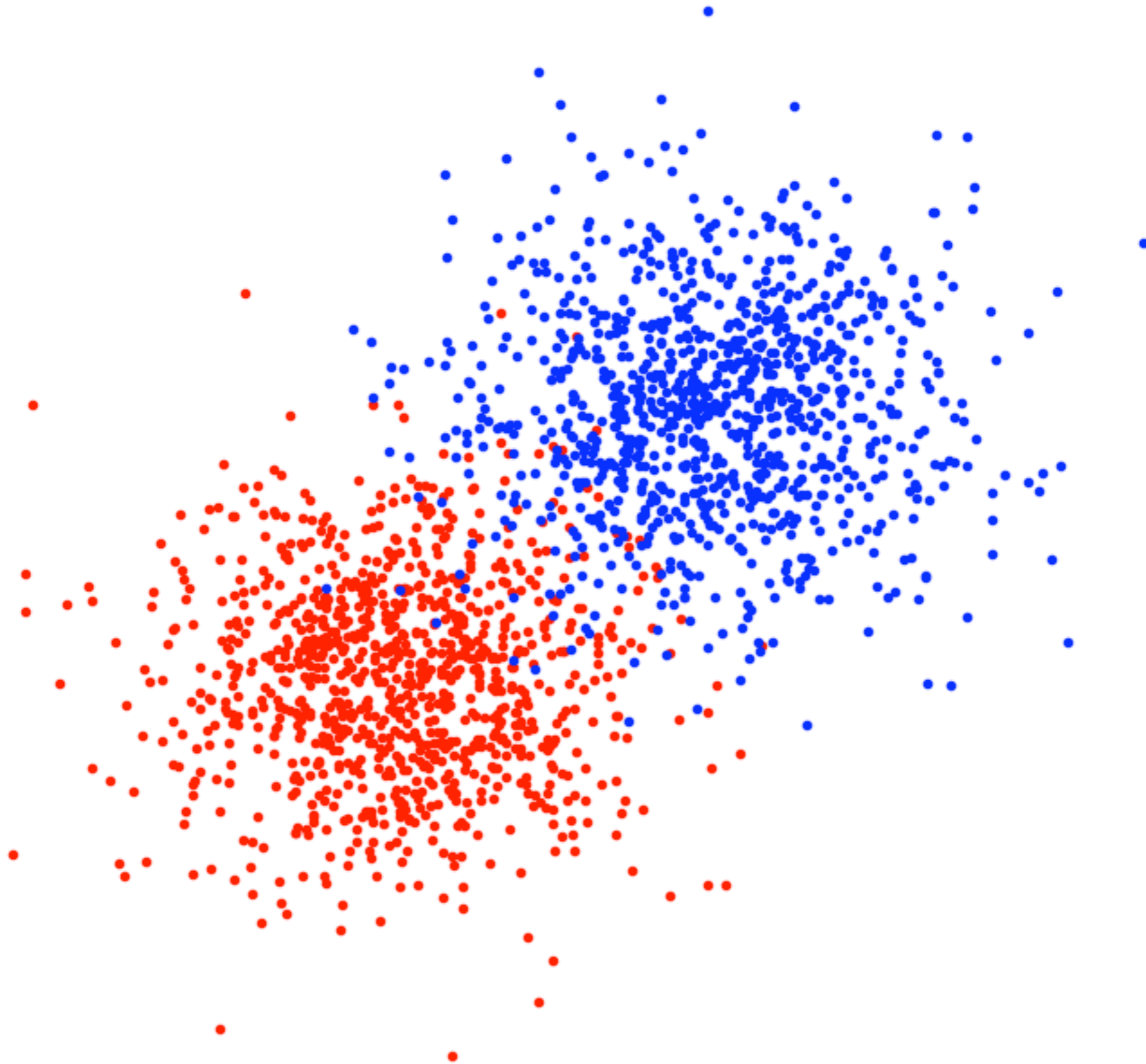- Is this greedy strategy **always** good?

# Extensions to Ponder

- What if labels are categorical but not binary?

- What if features are categorical but not binary?

- What if features are numeric?

- What if labels are numeric?

# Induction Learning

As you've seen, there are several issues that we must take into account when formalizing the notion of learning.

- The performance of the learning algorithm should be measured on unseen "test" data.

- The way in which we measure performance should depend on the problem we are trying to solve.

- There should be a strong relationship between the data that our algorithm sees at training time and the data it sees at test time.

# Setting up the learning problem

- Where does the data come from?

  - Data Generating Distribution

- How do we define that the model is good?

  - The loss function

- How do we combine those?

$$\epsilon \triangleq \mathbb{E}_{(x,y)\sim\mathcal{D}}\left[\ell(y, f(x))\right] = \sum_{(x,y)} \mathcal{D}(x,y)\ell(y, f(x))$$

So, putting it all together, we get a formal definition of induction machine learning: **Given (i) a loss function $\ell$ and (ii) a sample $D$ from some unknown distribution $\mathcal{D}$, you must compute a function $f$ that has low expected error $\epsilon$ over $\mathcal{D}$ with respect to $\ell$.**

So, putting it all together, we get a formal definition of induction machine learning: **Given (i) a loss function $\ell$ and (ii) a sample $D$ from some unknown distribution $\mathcal{D}$, you must compute a function $f$ that has low expected error $\epsilon$ over $\mathcal{D}$ with respect to $\ell$.**

**?** Why might it be a bad idea to use zero/one loss to measure performance for a regression problem?

So, putting it all together, we get a formal definition of induction machine learning: **Given (i) a loss function $\ell$ and (ii) a sample $D$ from some unknown distribution $\mathcal{D}$, you must compute a function $f$ that has low expected error $\epsilon$ over $\mathcal{D}$ with respect to $\ell$.**

?

Verify by calculation that we can write our training error as $\mathbb{E}_{(x,y)\sim D}\left[\ell(y, f(x))\right]$, by thinking of $D$ as a distribution that places probability $1/N$ to each example in $D$ and probability $0$ on everything else.