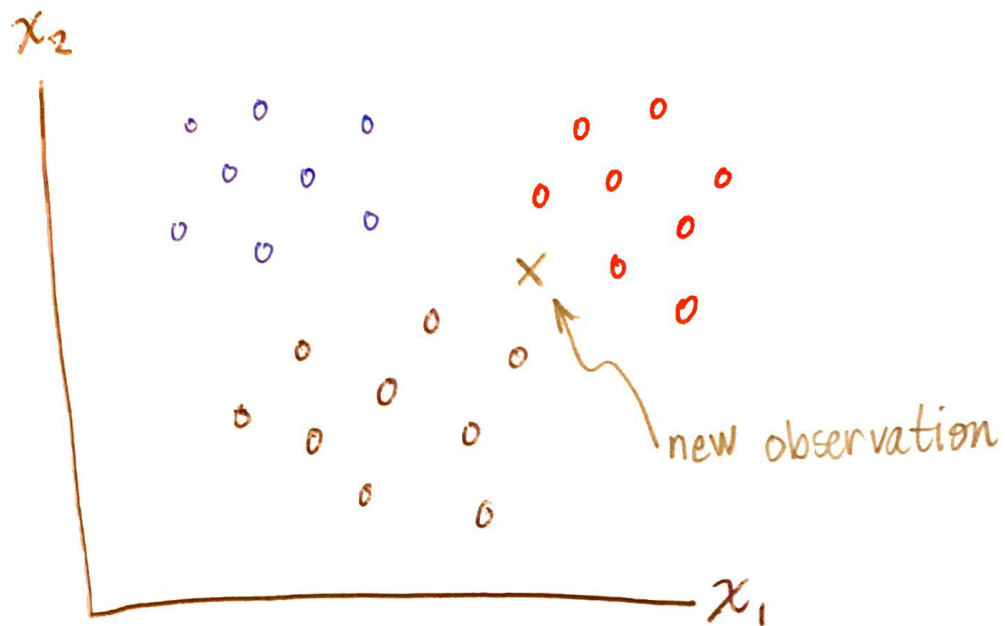


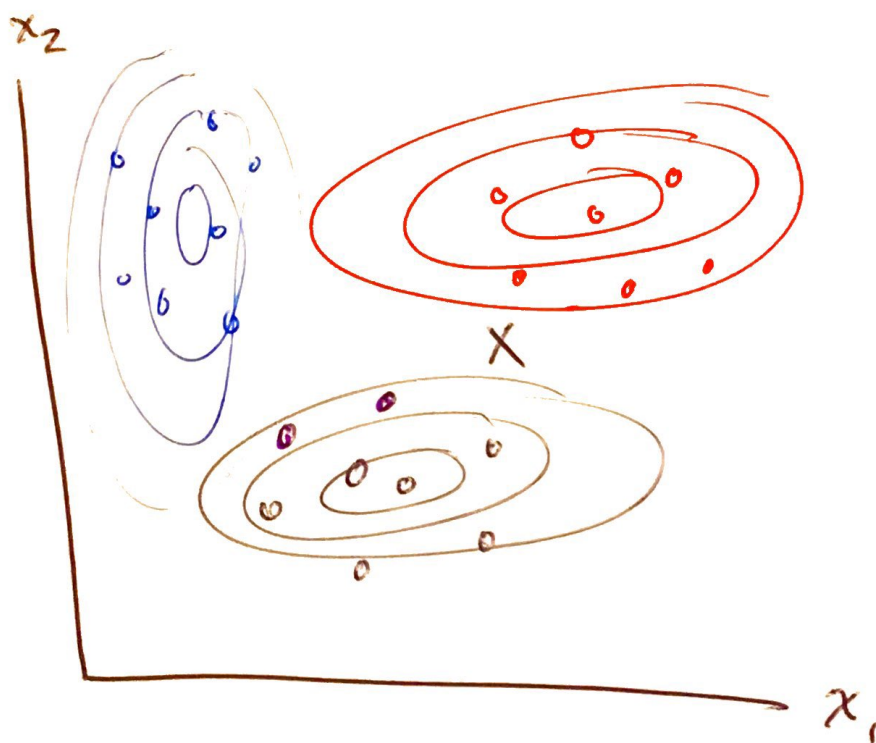
## Naive Bayes Classification



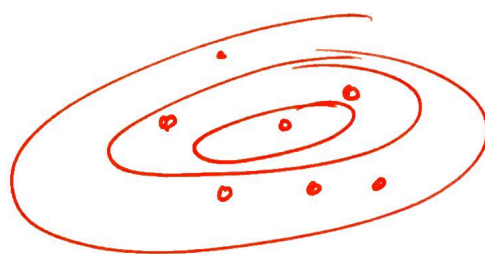
kNN and decision trees do not make any assumptions about the distribution of the underlying data  $\rightarrow$  which might help w/ outliers.

A naïve Bayes classifier assumes that  $x_1$  and  $x_2$  are independent, but it at least assumes some distribution on  $x_1$  and  $x_2$ .

Naive Bayes works by assuming that the observations are samples from distributions — a different distribution for each class.



If we knew those distributions,  
we could determine the probability  
of a new observation occurring at  
a specific location, ~~under~~ under the  
assumption that it is "red"

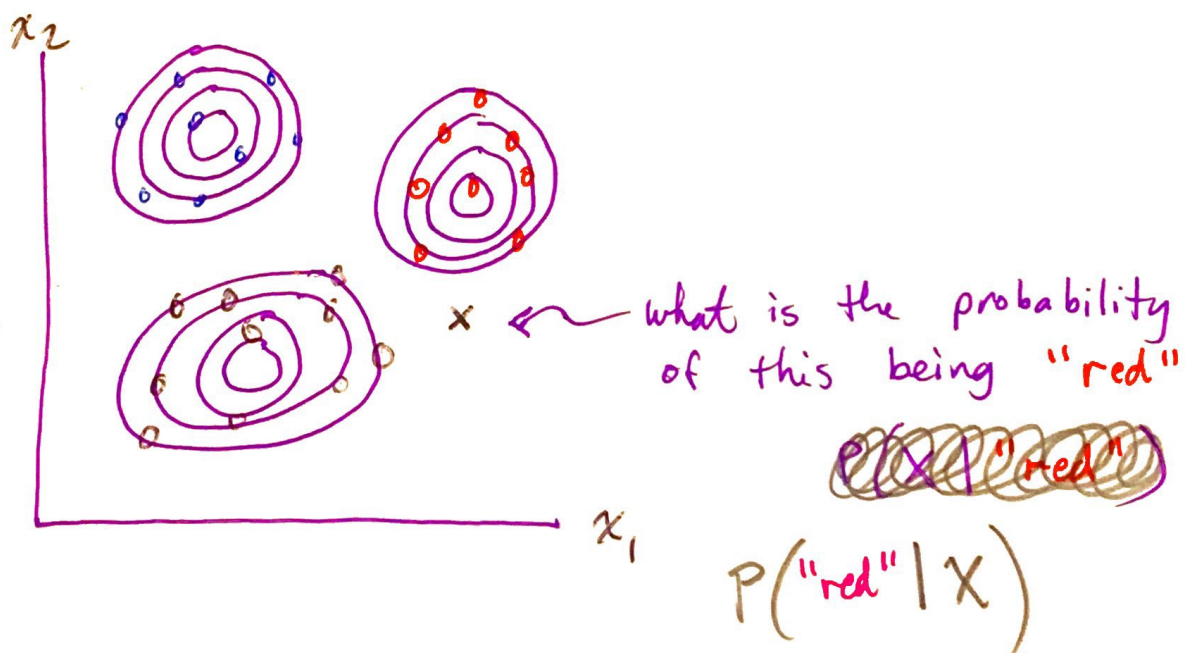


$X \leftarrow P(X | \text{"red"})$

what is the  
probability that  
it is here at  $X$   
given that it is  
"red" ?

-easy to calculate

Naive Bayes works by assuming  
that the observations are samples  
from distributions — a different  
distribution for each class.



vs.

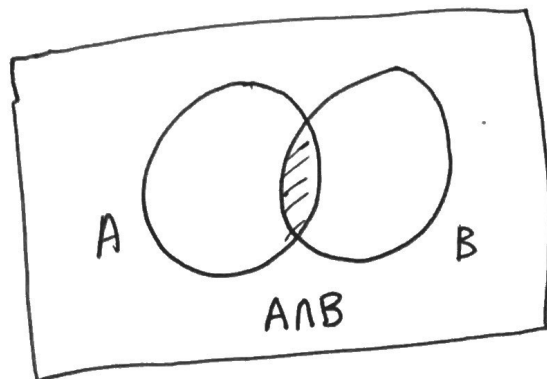
$$P(\text{"blue"} | x)$$

vs.

$$P(\text{"black"} | x)$$

## Bayes Rule

to the rescue!



$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (*)$$

Likewise

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad (**)$$

Solve  $(**)$  for  $P(A \cap B)$

$$P(A \cap B) = P(B|A) \cdot P(A)$$

Subst. into  $(*)$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

↑  
what we want

↙ what we know

Construct a naive Bayes classifier, using `fitcnb`

```
>> mdl = fitcnb(dataTrain, "group")
```

Recall setup:

```
groupData = readtable("groupData.csv")  
groupData.group = categorical(groupData.group)  
pt = cvpartition(groupData.group, "Holdout", 0.35)  
dataTrain = groupData(training(pt), :)  
dataTest = groupData(test(pt), :)
```

By default, the numeric predictors are modeled as normal distributions. Try changing the distribution to "kernel"

```
>> mdl = fitcnb(dataTrain, "group", ...  
    "DistributionNames", "normal")
```

Plot the results

```
>> label = predict (mdl, dataTest)
```

```
>> gscatter (groupData.x, groupData.y, ...  
            groupData.group)
```

% now plot the predicted labels on top  
of the original data

```
>> hold on
```

```
>> gscatter (dataTest.x, dataTest.y, ...
```

label, [], "0", 15)

plot  
predicted  
labels

default  
colors

little  
0's

marker  
size

Calculate error (loss)

~~predGroups = predict(mdl, dataTest)~~

>> err = loss(mdl, dataTest)

| Classifier              | Loss   |
|-------------------------|--------|
| Naive Bayes<br>"kernel" | 0.1254 |
| Naive Bayes<br>"normal" | 0.1593 |
| Tree                    | 0.1259 |
| kNN<br>"k=10"           | 0.1171 |



## Heart Disease Analysis

Read and load data:

```
>> heartData = readTable("heartDataNum.csv")  
>> heartData, HeartDisease = categorical(...  
    heartData, HeartDisease)
```

Partition into training and test sets:

```
>> pt = cvpartition(heartData, HeartDisease, ...  
    "HoldOut", 0.3)  
>> hdTrain = heartData(training(pt), :)  
>> hdTest = heartData(test(pt), :)
```

Create a naive Bayes classification model

```
>> mdl = fitcnb(hdTrain, "HeartDisease")
```

Calculate error

% training error

err Train = resubLoss (mdl)

% test error

err Test = loss (mdl, hd Test)

## Heart Disease - Numeric and Categorical

```
>> heartData = readtable("heartDataAll.csv")
```

```
% convert text labels to categorical arrays
```

```
>> heartData = convertvars(heartData, ...  
    12:22, "categorical")
```

```
% partition the data into training and test  
>> pt = cvpartition(heartData.HeartDisease, "HoldOut", 0.3)
```

```
>> hdTrain = heartData(training(pt), :)
```

```
>> hdTest = heartData(test(pt), :)
```

```
% construct a naïve Bayes model on the  
    categorical & numeric features
```

```
mdl = fitcnb(hdTrain, "HeartDisease")
```

% Use "kernel" for the numeric variables  
(the first 11)

and "mvnm" multivariate, multinomial  
for the categorical variables  
(the last 10)

```
>> dists = [ repmat("kernel", 1, 11)  
             repmat("mvnm", 1, 10) ]
```

dists = 1x21 string  
"kernel" "kernel" ... "mvnm" "mvnm"

```
>> mdl = fitnb(hdTrain, "Heart Disease", ...  
               "Distribution Names", dists)
```

% training and testing error

```
>> errTrain = resubLoss(mdl)
```

```
>> errTest = loss(mdl, hdTest)
```