

Naive Bayes



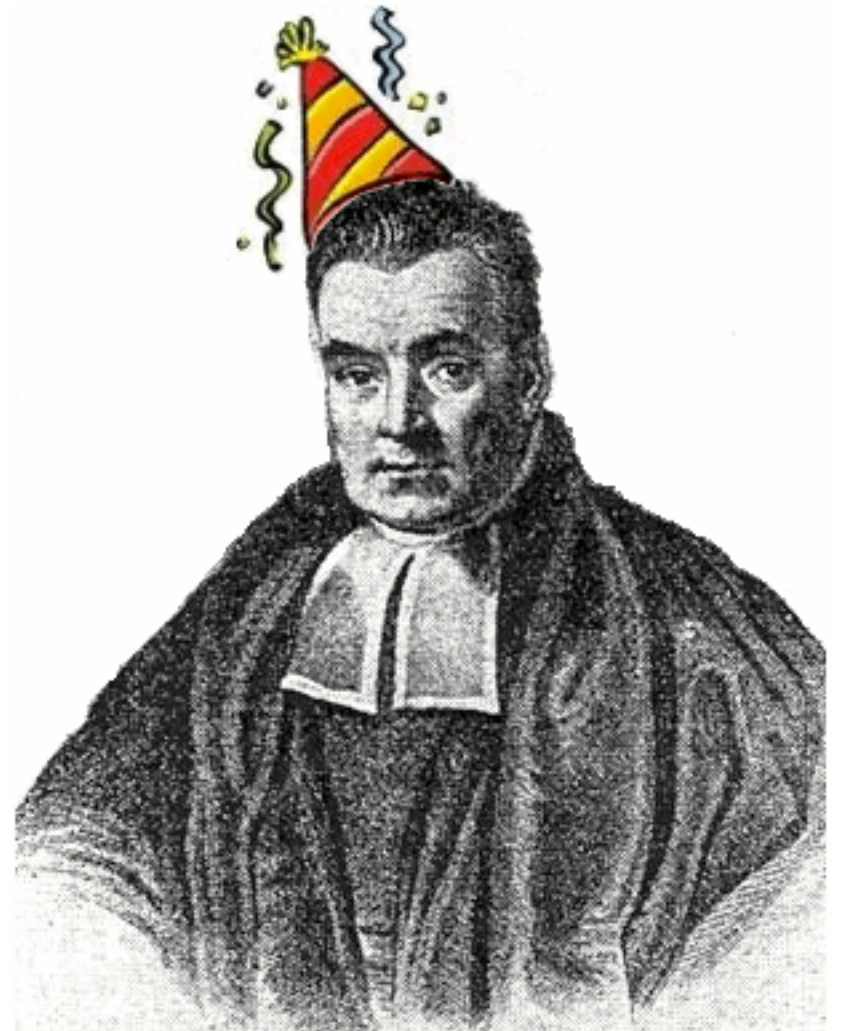
Naïve



Updating the state of
knowledge

step by step

with new information



What is classification?

Deciding among hypotheses (labels),
using information we have (features)
for each example

3 Features: Votes on 3 Bills

2 Labels: Democrat / Republican

Prediction:

I know your votes, I'm trying to guess your party

3 Features: Votes on 3 Bills

2 Labels: Democrat / Republican

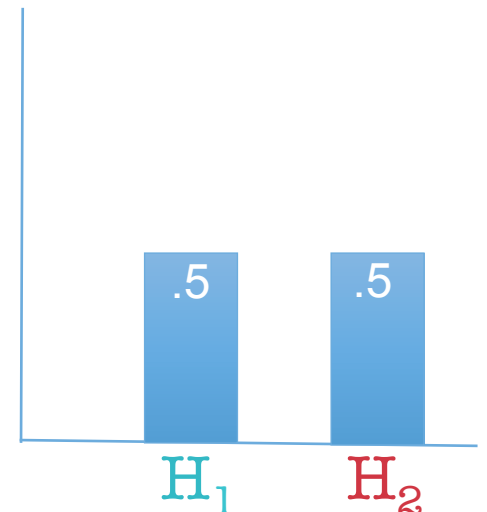
Prediction:

I know your votes, I'm trying to guess your party

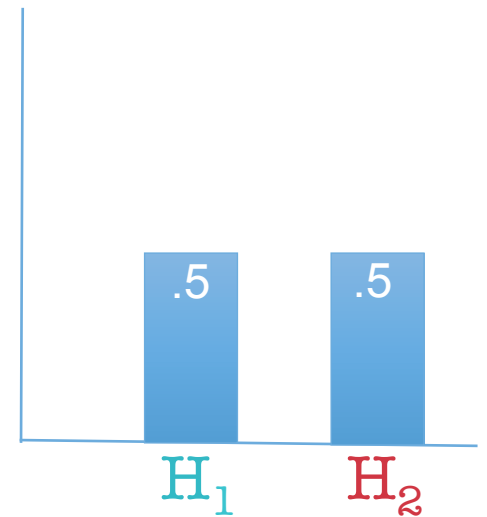
2 Labels

H_1 : Democrat

H_2 : Republican



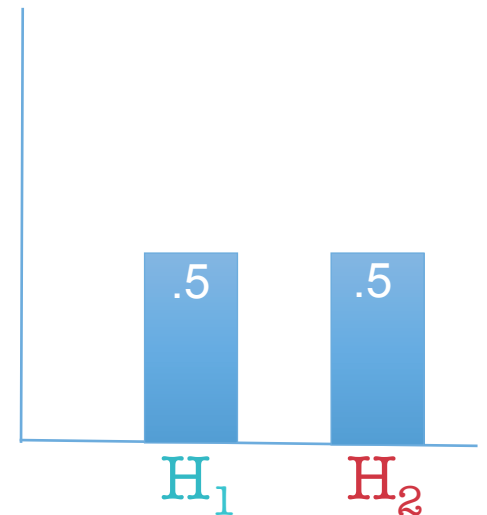
Prior: Initial belief



Prior: Initial belief

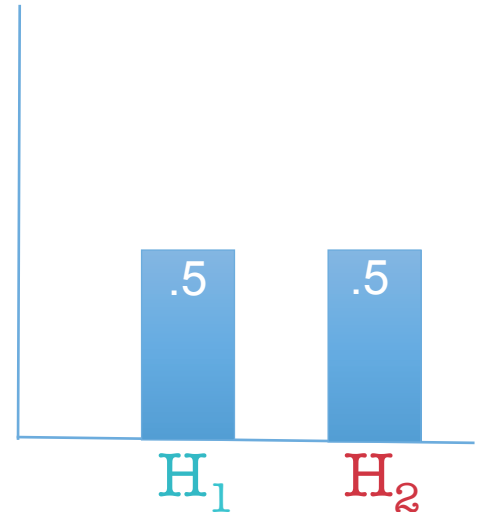
50 - 50 ? $P(\text{Democrat}) = 0.5$?

(Uninformative prior)



Prior: Initial belief

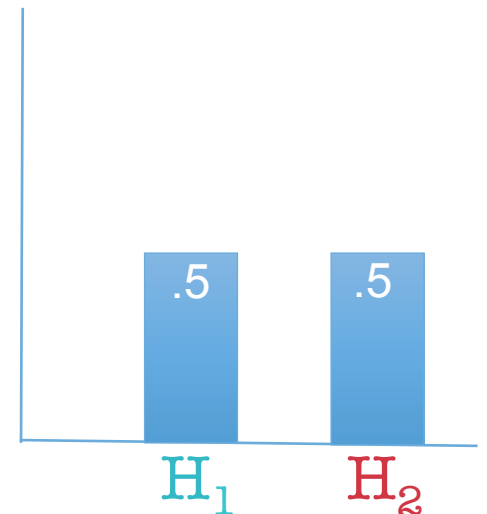
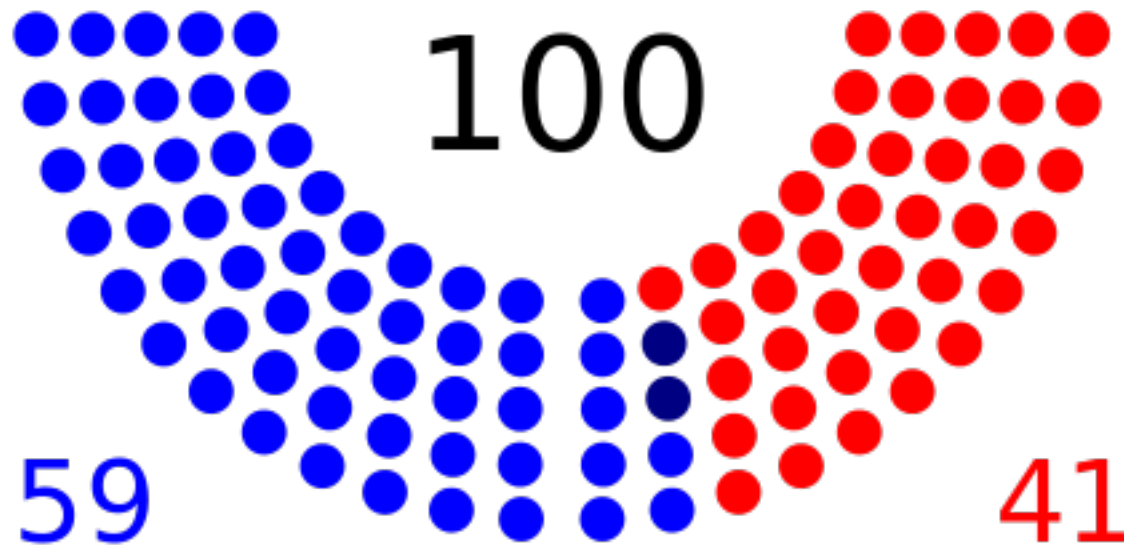
What's my best guess without any vote info?



Prior: Initial belief

What's my best guess without any vote info?

I'd guess **democrat** since there are more of them.

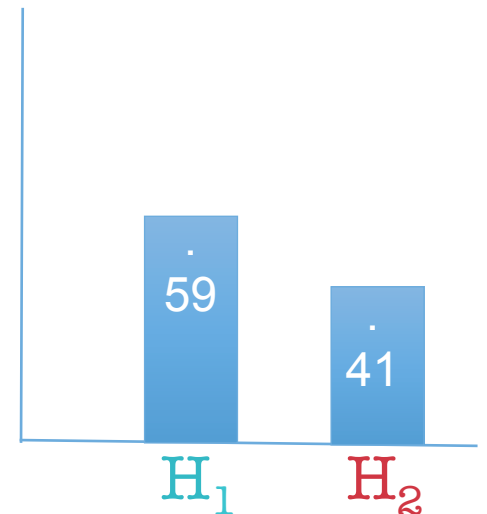
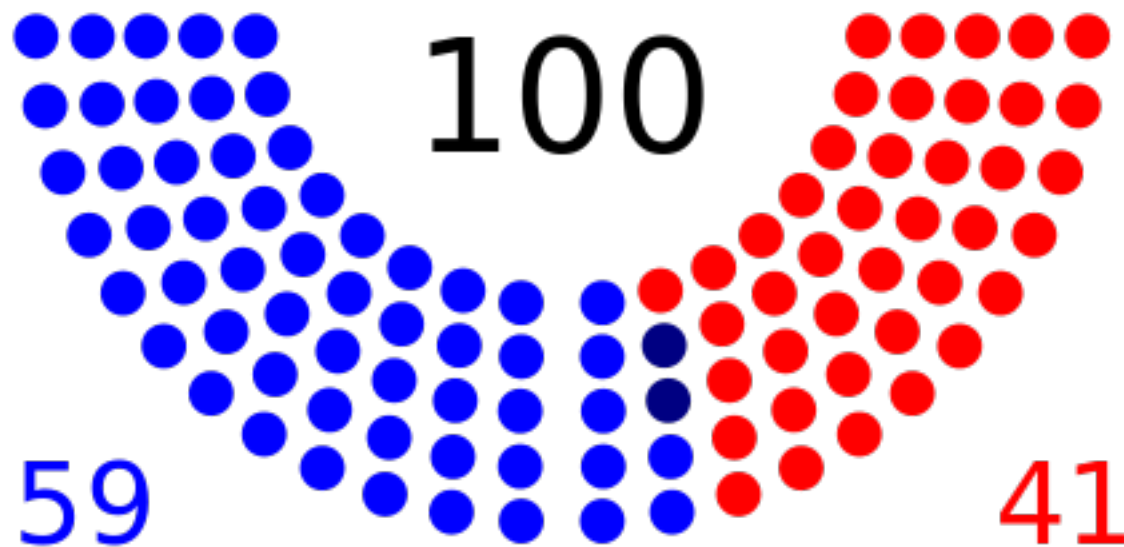


Prior: Initial belief

What's my best guess without any vote info?

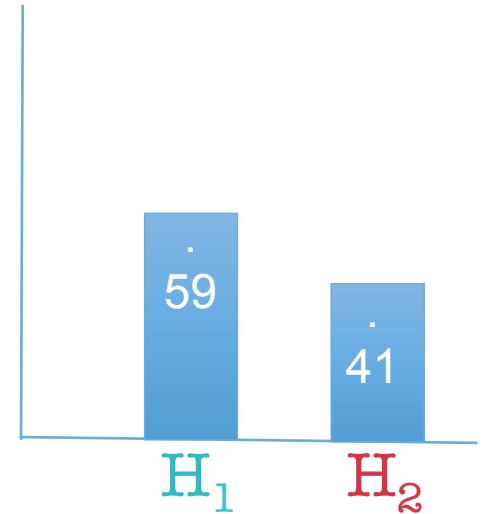
I'd guess democrat since there are more of them.

$$P(\text{Democrat}) = 0.59$$



Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

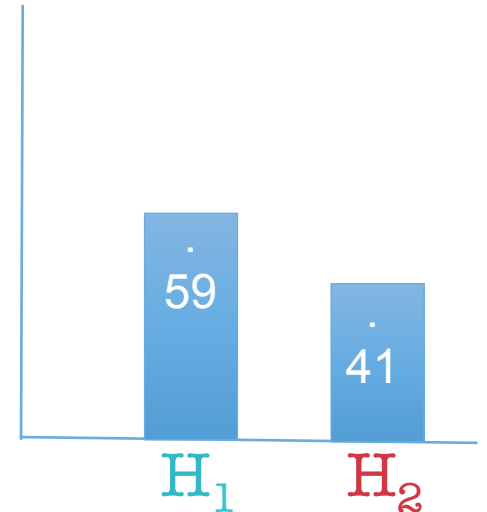


Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

Voted YES on Net Neutrality



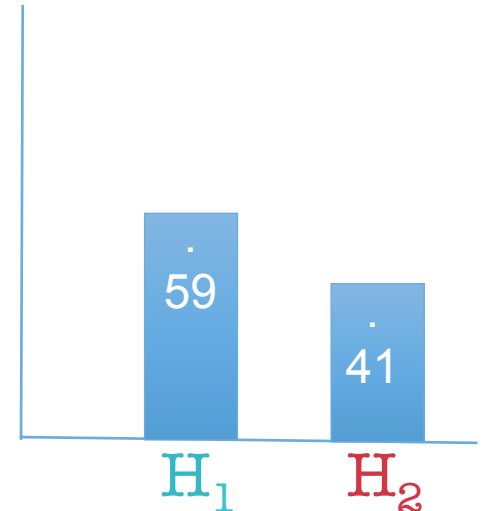
Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{\text{NN}}) = \frac{P(Y_{\text{NN}} | \text{Dem}) P(\text{Dem})}{P(Y_{\text{NN}})}$$



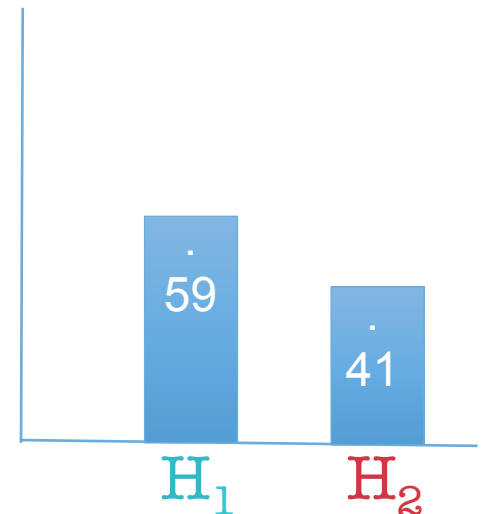
Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{\text{NN}}) = \frac{\overset{\text{likelihood}}{P(Y_{\text{NN}} | \text{Dem})} \overset{\text{prior}}{P(\text{Dem})}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{\text{NN}})}}$$



Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

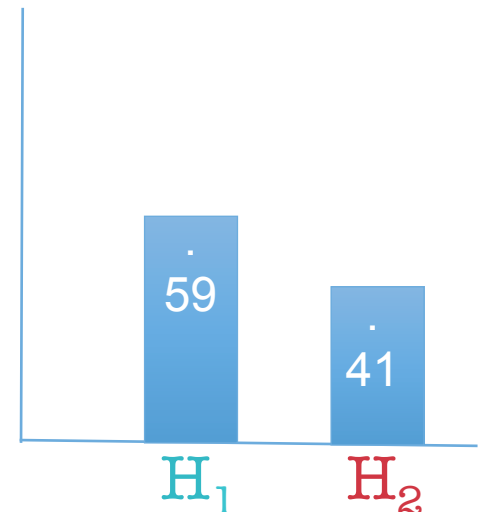
Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{NN}) = \frac{\overset{\text{likelihood}}{P(Y_{NN} | \text{Dem})} \overset{\text{prior}}{P(\text{Dem})}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{NN})}}$$

posterior

$$P(Y_{NN} | \text{Dem})$$

Prob. of voting yes
on net neutrality
if you're democrat



Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

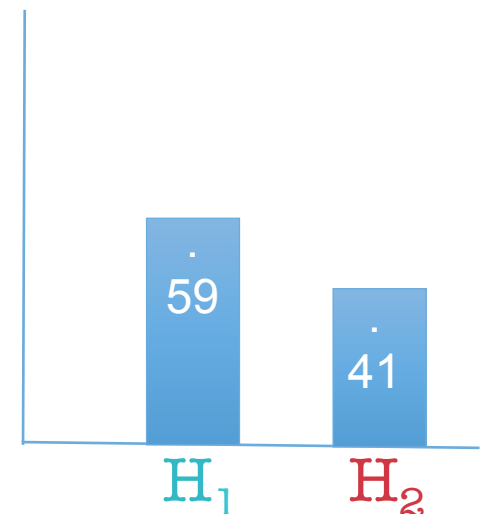
Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{NN}) = \frac{\overset{\text{posterior}}{P(\text{Dem} | Y_{NN})} = \frac{\overset{\text{likelihood}}{P(Y_{NN} | \text{Dem})} \overset{\text{prior}}{P(\text{Dem})}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{NN})}}$$

$$P(\text{Rep} | Y_{NN}) = \frac{P(Y_{NN} | \text{Rep}) P(\text{Rep})}{P(Y_{NN})}$$

$$P(Y_{NN} | \text{Dem})$$

Prob. of voting yes
on net neutrality
if you're democrat



Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

Voted YES on Net Neutrality

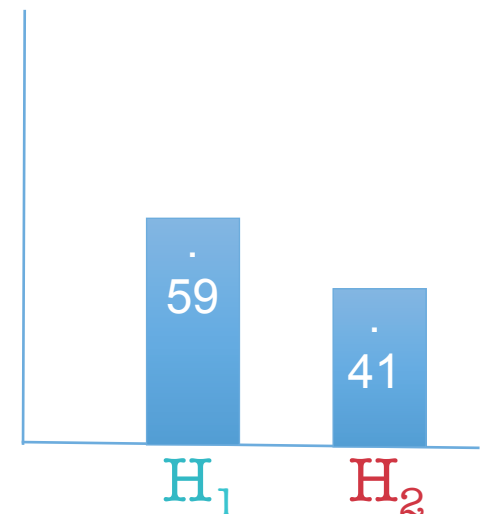
$$P(\text{Dem} | Y_{NN}) = \frac{\overset{\text{posterior}}{P(\text{Dem} | Y_{NN})} = \frac{\overset{\text{likelihood}}{P(Y_{NN} | \text{Dem})} \overset{\text{prior}}{P(\text{Dem})}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{NN})}}$$

$$P(\text{Rep} | Y_{NN}) = \frac{P(Y_{NN} | \text{Rep}) P(\text{Rep})}{P(Y_{NN})}$$

$$P(Y_{NN} | \text{Dem})$$

Prob. of voting yes
on net neutrality
if you're democrat

$$P(Y_{NN} | \text{Rep})$$

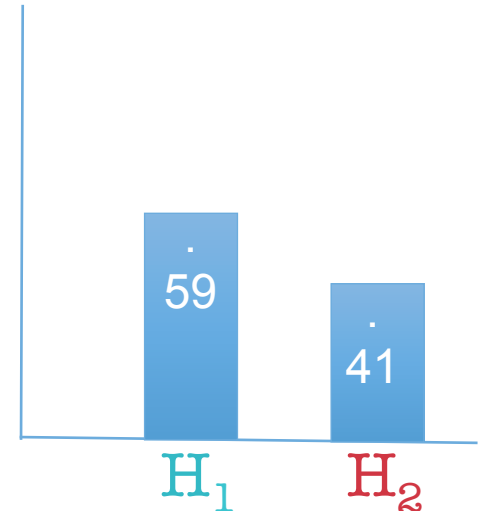


$$P(Y_{NN}|\text{Dem})$$

Prob. of voting yes
on net neutrality
if you're democrat

$$P(Y_{NN}|\text{Rep})$$

Prob. of voting
yes
on net neutrality
if you're
republican



Training set has the answers!

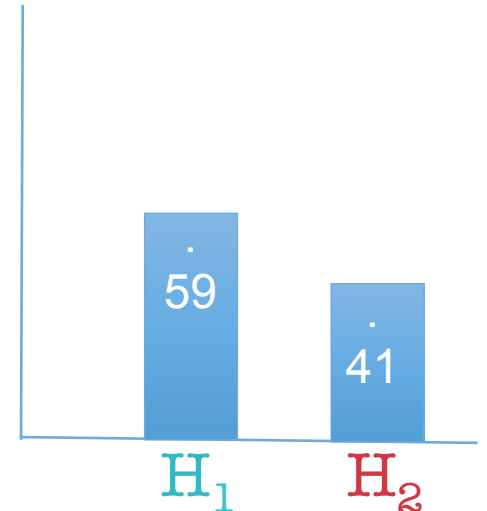
We know Dem/Rep for each person, we know their votes!

$$P(Y_{NN}|\text{Dem})$$

Prob. of voting yes
on net neutrality
if you're democrat

$$P(Y_{NN}|\text{Rep})$$

Prob. of voting
yes
on net neutrality
if you're
republican



Training set has the answers!

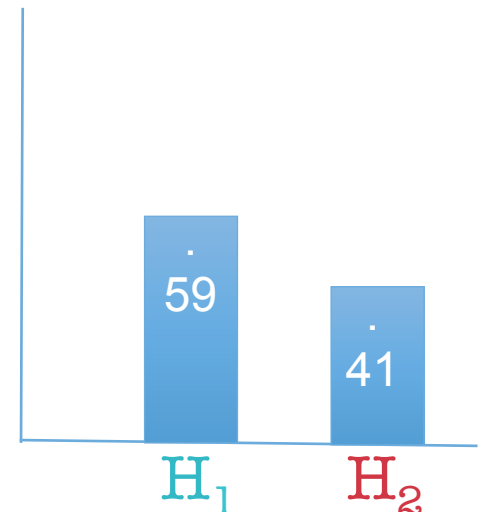
We know Dem/Rep for each person, we know their votes!

$$P(Y_{NN} | \text{Dem})$$

$$\approx \frac{\text{\# democrats that } Y_{NN}}{\text{\# all democrats}}$$

$$P(Y_{NN} | \text{Rep})$$

$$\approx \frac{\text{\# republicans that } Y_{NN}}{\text{\# all republicans}}$$



Training set has the answers!

We know Dem/Rep for each person, we know their votes!

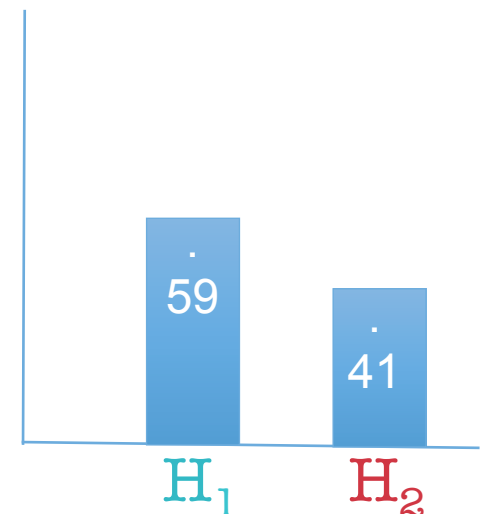
$$P(Y_{NN} | \text{Dem})$$

$$\approx \frac{\text{\# democrats that } Y_{NN}}{\text{\# all democrats}}$$

$$P(Y_{NN} | \text{Rep})$$

$$\approx \frac{\text{\# republicans that } Y_{NN}}{\text{\# all republicans}}$$

For likelihoods of discrete data,
training/fitting means counting!
(and estimating likelihoods by dividing counts)



Training set has the answers!

We know Dem/Rep for each person, we know their votes!

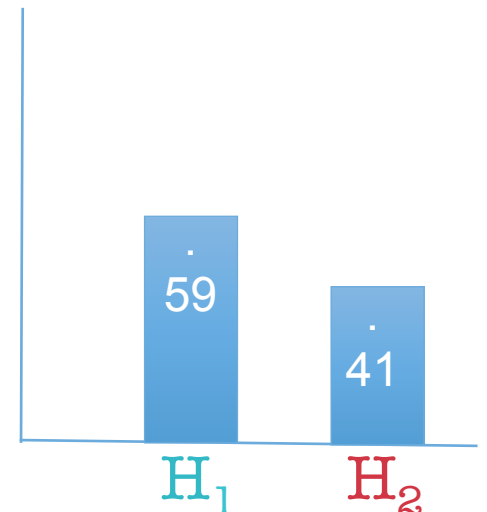
$$P(Y_{NN}|\text{Dem})$$

$$\approx \frac{56}{59} = 0.949$$

$$P(Y_{NN}|\text{Rep})$$

$$\approx \frac{34}{41} = 0.829$$

For likelihoods of discrete data,
training/fitting means counting!
(and estimating likelihoods by dividing counts)



Prior: Initial belief

$$P(\text{Democrat}) = 0.59$$

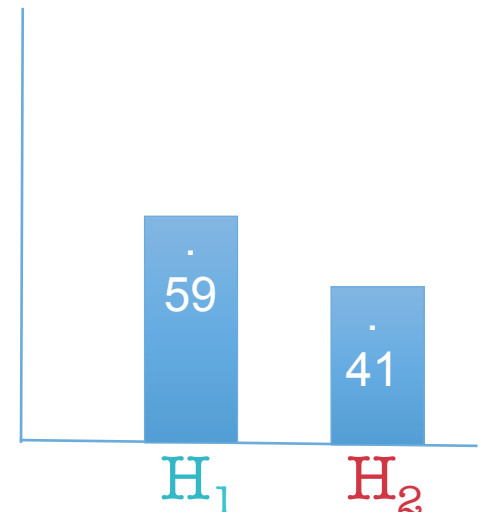
New information (feature 1):

Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{\text{NN}}) = \frac{\overset{\text{likelihood}}{P(Y_{\text{NN}} | \text{Dem})} \overset{\text{prior}}{P(\text{Dem})}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{\text{NN}})}}$$

evidence
(normalization factor)

$$P(\text{Rep} | Y_{\text{NN}}) = \frac{P(Y_{\text{NN}} | \text{Rep}) P(\text{Rep})}{P(Y_{\text{NN}})}$$



Prior: Initial belief

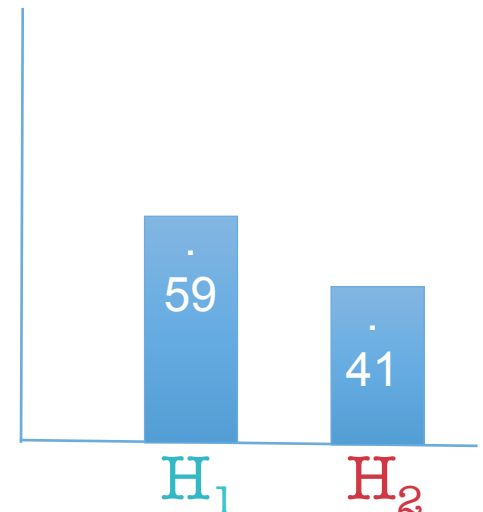
$$P(\text{Democrat}) = 0.59$$

New information (feature 1):

Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{\text{NN}}) = \frac{\overset{\text{likelihood}}{0.949} * \overset{\text{prior}}{0.59}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{\text{NN}})}}$$

$$P(\text{Rep} | Y_{\text{NN}}) = \frac{0.829 * 0.41}{P(Y_{\text{NN}})}$$



Current belief

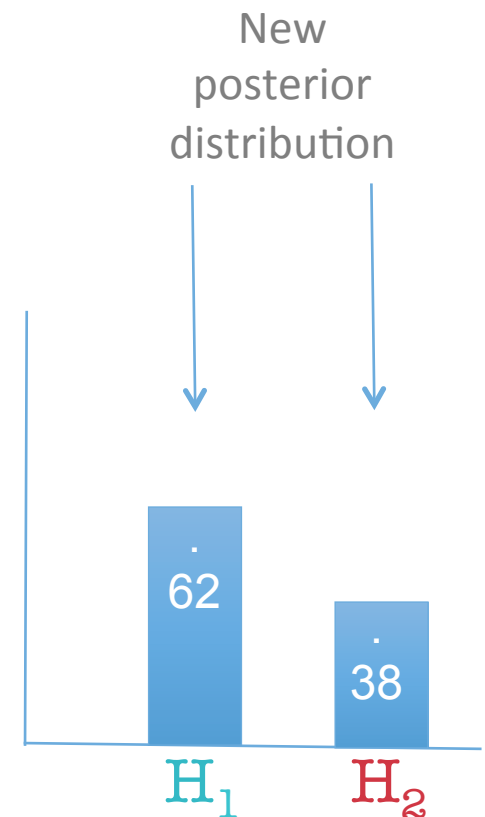
$$P(\text{Democrat} | Y_{NN}) = 0.62$$

New information (feature 1):

Voted YES on Net Neutrality

$$P(\text{Dem} | Y_{NN}) = \frac{\overset{\text{likelihood}}{0.949} * \overset{\text{prior}}{0.59}}{\underset{\substack{\text{evidence} \\ \text{(normalization factor)}}}{P(Y_{NN})}}$$

$$P(\text{Rep} | Y_{NN}) = \frac{0.829 * 0.41}{P(Y_{NN})}$$

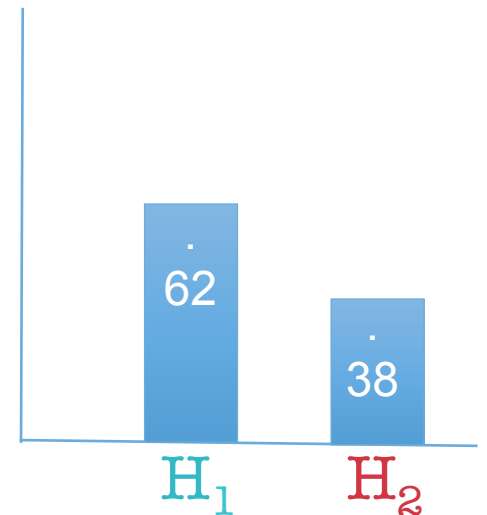


Current belief

$$P(\text{Democrat} | Y_{NN}) = 0.62$$

New information (feature 2):

Voted YES on Tax Cuts



Current belief

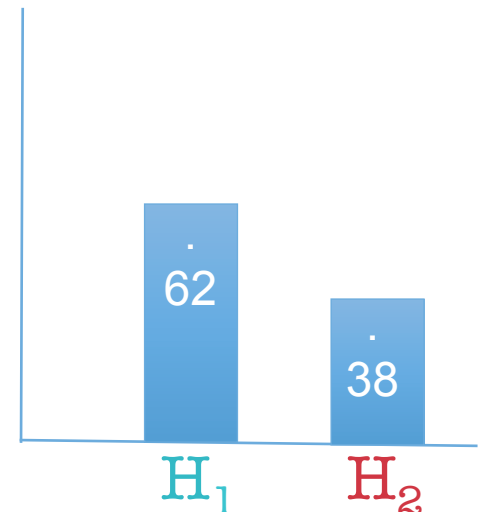
$$P(\text{Democrat} | Y_{NN}) = 0.62$$

New information (feature 2):

Voted YES on Tax Cuts

$$P(\text{Dem} | Y_{NN}, Y_{TC}) = \frac{P(Y_{TC} | \text{Dem}) P(\text{Dem} | Y_{NN})}{P(Y_{TC})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}) = \frac{P(Y_{TC} | \text{Rep}) P(\text{Rep} | Y_{NN})}{P(Y_{TC})}$$



Current belief

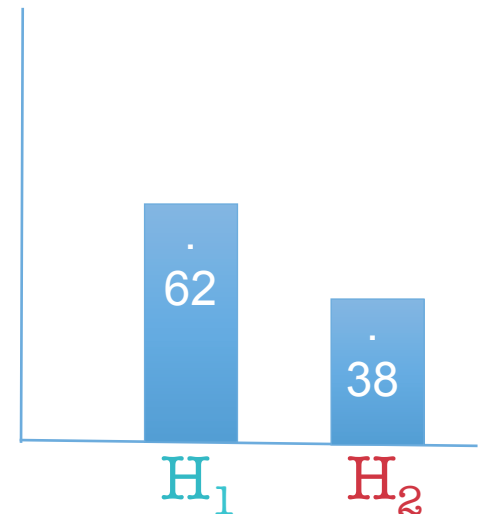
$$P(\text{Democrat} | Y_{NN}) = 0.62$$

$$P(Y_{TC} | \text{Dem})$$

$$\approx \frac{10}{59} = 0.169$$

$$P(Y_{TC} | \text{Rep})$$

$$\approx \frac{35}{41} = 0.854$$



Current belief

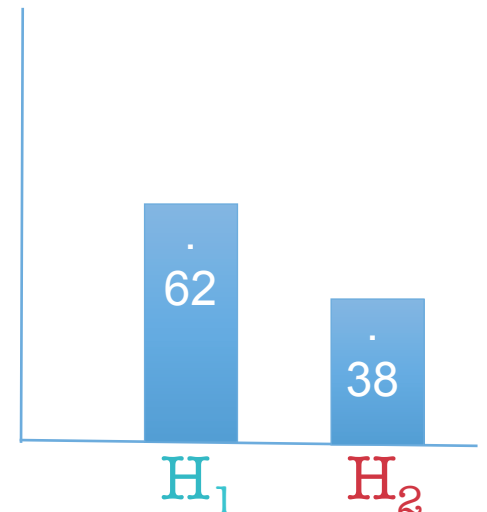
$$P(\text{Democrat} | Y_{NN}) = 0.62$$

New information (feature 2):

Voted YES on Tax Cuts

$$P(\text{Dem} | Y_{NN}, Y_{TC}) = \frac{P(Y_{TC} | \text{Dem}) P(\text{Dem} | Y_{NN})}{P(Y_{TC})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}) = \frac{P(Y_{TC} | \text{Rep}) P(\text{Rep} | Y_{NN})}{P(Y_{TC})}$$



Current belief

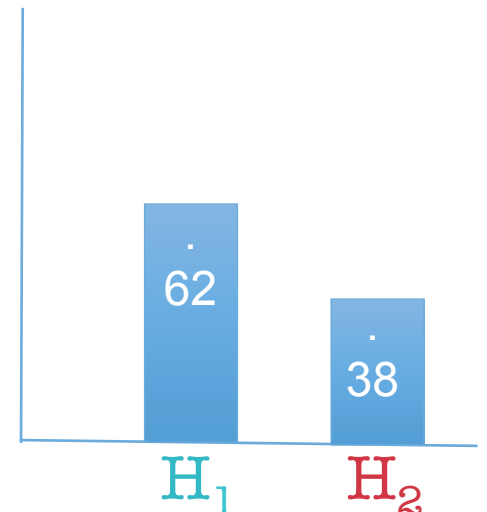
$$P(\text{Democrat} | Y_{NN}) = 0.62$$

New information (feature 2):

Voted YES on Tax Cuts

$$P(\text{Dem} | Y_{NN}, Y_{TC}) = \frac{0.169 * 0.62}{P(Y_{TC})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}) = \frac{0.854 * 0.38}{P(Y_{TC})}$$



Current belief

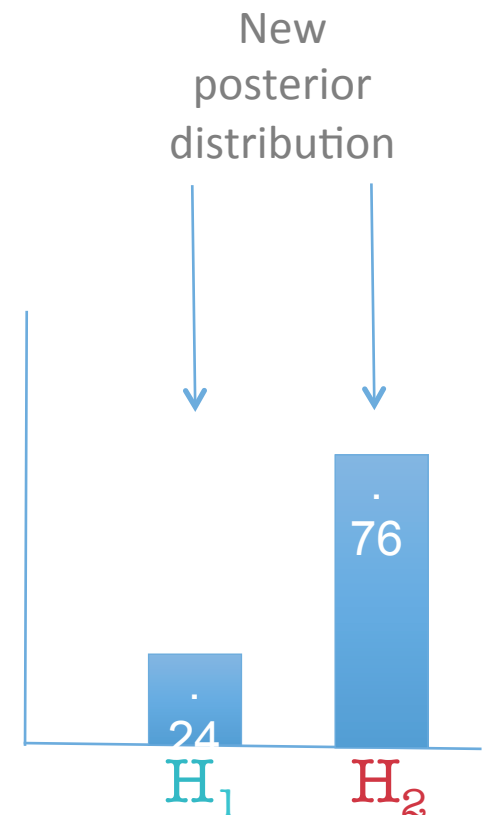
$$P(\text{Democrat} | Y_{NN}, Y_{TC}) = 0.24$$

New information (feature 2):

Voted YES on Tax Cuts

$$P(\text{Dem} | Y_{NN}, Y_{TC}) = \frac{0.169 * 0.62}{P(Y_{TC})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}) = \frac{0.854 * 0.38}{P(Y_{TC})}$$

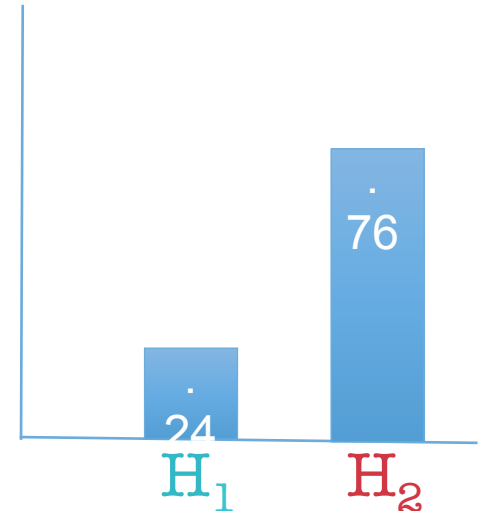


Current belief

$$P(\text{Democrat} | Y_{\text{NN}}, Y_{\text{TC}}) = 0.24$$

New information (feature 3):

Voted NO on License-free Guns



Current belief

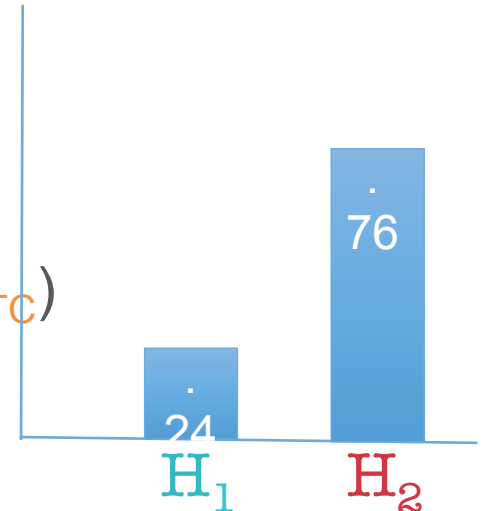
$$P(\text{Democrat} | Y_{NN}, Y_{TC}) = 0.24$$

New information (feature 3):

Voted NO on License-free Guns

$$P(\text{Dem} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{P(N_{LG} | \text{Dem}) P(\text{Dem} | Y_{NN}, Y_{TC})}{P(N_{LG})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{P(N_{LG} | \text{Rep}) P(\text{Rep} | Y_{NN}, Y_{TC})}{P(N_{LG})}$$



Current belief

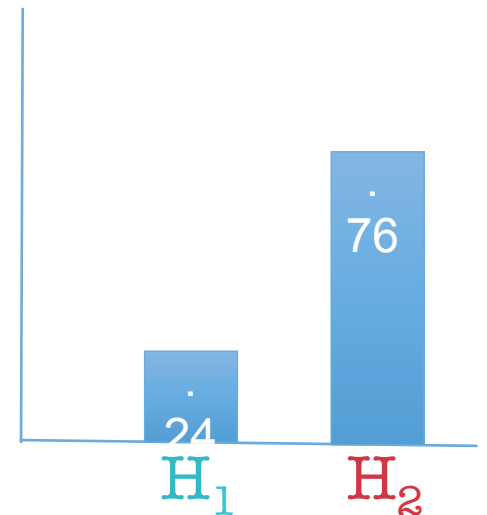
$$P(\text{Democrat} | Y_{\text{NN}}, Y_{\text{TC}}) = 0.24$$

$$P(N_{\text{LG}} | \text{Dem})$$

$$\approx \frac{53}{59} = 0.898$$

$$P(N_{\text{LG}} | \text{Rep})$$

$$\approx \frac{23}{41} = 0.561$$



Current belief

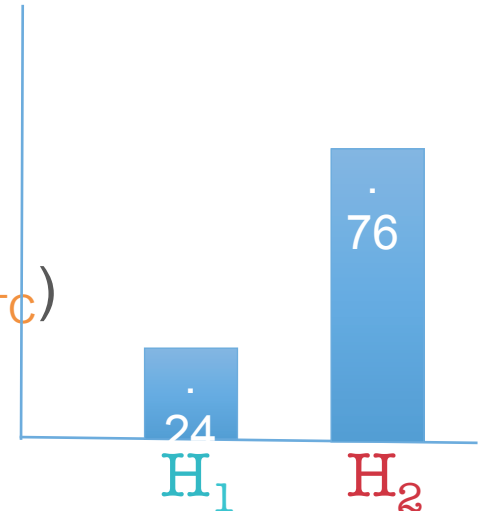
$$P(\text{Democrat} | Y_{NN}, Y_{TC}) = 0.24$$

New information (feature 3):

Voted NO on License-free Guns

$$P(\text{Dem} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{P(N_{LG} | \text{Dem}) P(\text{Dem} | Y_{NN}, Y_{TC})}{P(N_{LG})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{P(N_{LG} | \text{Rep}) P(\text{Rep} | Y_{NN}, Y_{TC})}{P(N_{LG})}$$



Current belief

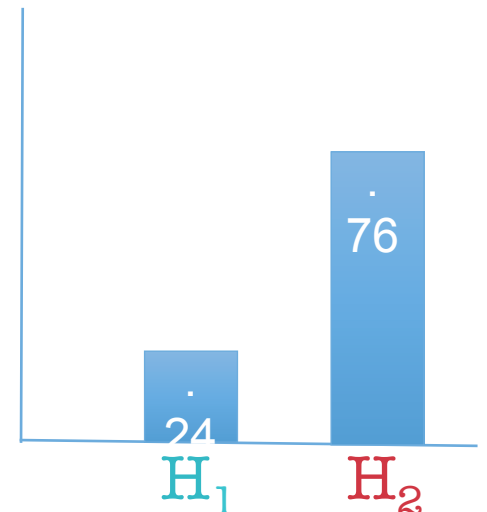
$$P(\text{Democrat} | Y_{NN}, Y_{TC}) = 0.24$$

New information (feature 3):

Voted NO on License-free Guns

$$P(\text{Dem} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{0.898 * 0.24}{P(N_{LG})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{0.561 * 0.76}{P(N_{LG})}$$



Current belief

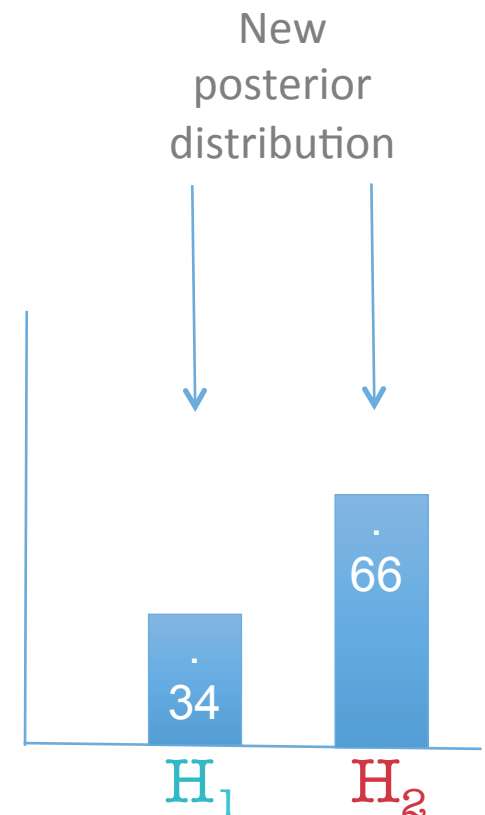
$$P(\text{Democrat} | Y_{NN}, Y_{TC}, N_{LG}) = 0.34$$

New information (feature 3):

Voted NO on License-free Guns

$$P(\text{Dem} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{0.898 * 0.24}{P(N_{LG})}$$

$$P(\text{Rep} | Y_{NN}, Y_{TC}, N_{LG}) = \frac{0.561 * 0.76}{P(N_{LG})}$$



Current belief

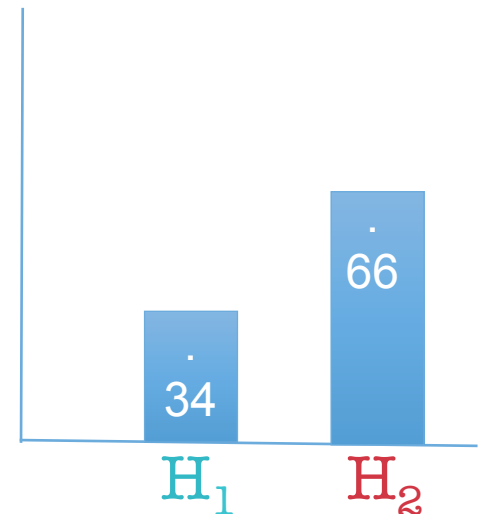
$$P(\text{Democrat} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = 0.34$$

Classify this person that voted

Yes on Net Neutrality (Y_{NN}),

Yes on Tax Cuts (Y_{TC}),

No on License-free Guns (N_{LG})



Current belief

$$P(\text{Democrat} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = 0.34$$

Classify this person that voted

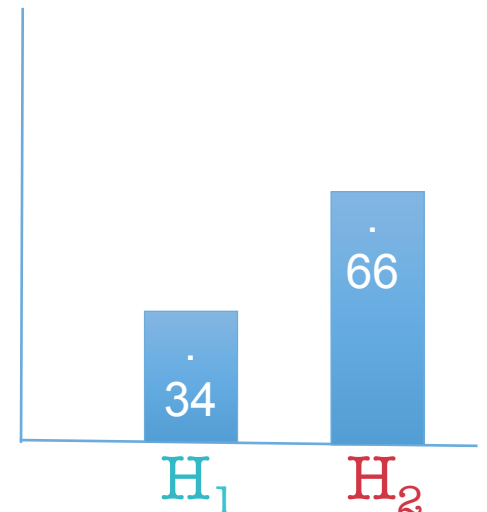
Yes on Net Neutrality (Y_{NN}),

Yes on Tax Cuts (Y_{TC}),

No on License-free Guns (N_{LG})

My strongest belief is in H_2 ,

I classify this person with the label **Republican**.



Naïve Bayes

Training:

Count and calculate the likelihood of each feature value for each class:

$$P(Y_{NN}|\text{Dem}) = 1 - P(N_{NN}|\text{Dem})$$

$$P(Y_{NN}|\text{Rep}) = 1 - P(N_{NN}|\text{Rep})$$

$$P(Y_{TC}|\text{Dem}) = 1 - P(N_{TC}|\text{Dem})$$

$$P(Y_{TC}|\text{Rep}) = 1 - P(N_{TC}|\text{Rep})$$

$$P(Y_{LG}|\text{Dem}) = 1 - P(N_{LG}|\text{Dem})$$

$$P(Y_{LG}|\text{Rep}) = 1 - P(N_{LG}|\text{Rep})$$

Prediction:

Use Bayes to update priors with the likelihoods,
Pick label with the highest posterior probability.

What was the naïve part?



We assumed each feature is independent

We assumed each feature is independent

Easier to see in a single update rather than
sequential

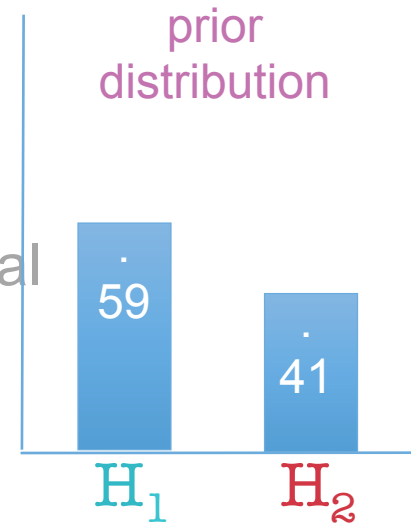
We assumed each feature is independent

Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG

$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) P(\text{Dem})}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

We assumed each feature is independent

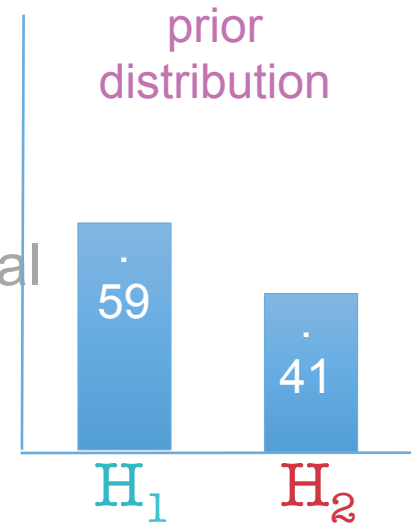
Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG



$$\text{posterior } P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{\text{likelihood } P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) \text{ prior } P(\text{Dem})}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

We assumed each feature is independent

Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG



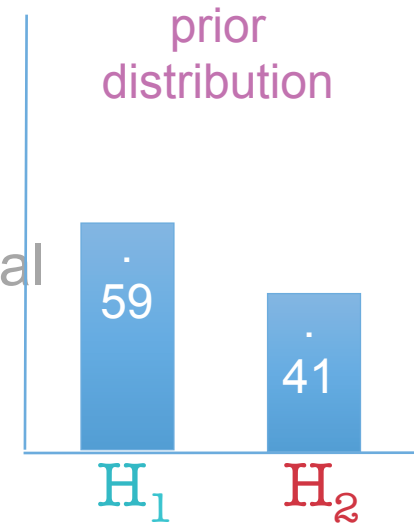
$$\text{posterior } P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{\text{likelihood } P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) \text{ prior } P(\text{Dem})}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

Independence Assumption:

$$P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) = P(Y_{\text{NN}} | \text{Dem}) P(Y_{\text{TC}} | \text{Dem}) P(N_{\text{LG}} | \text{Dem})$$

We assumed each feature is independent

Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG



$$\text{posterior } P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{\text{likelihood } P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) \text{ prior } P(\text{Dem})}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

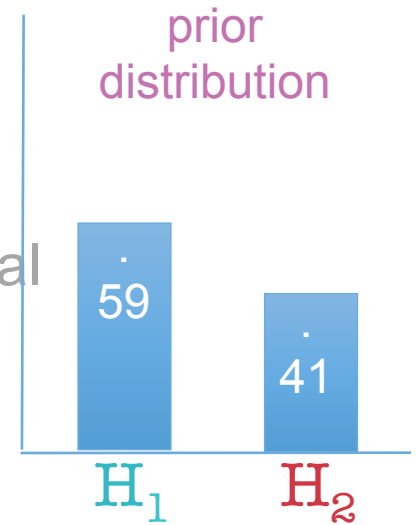
Independence Assumption:

$$P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}} | \text{Dem}) = P(Y_{\text{NN}} | \text{Dem}) P(Y_{\text{TC}} | \text{Dem}) P(N_{\text{LG}} | \text{Dem})$$

Not even close in most cases!
Naïve Bayes still works well.

We assumed each feature is independent

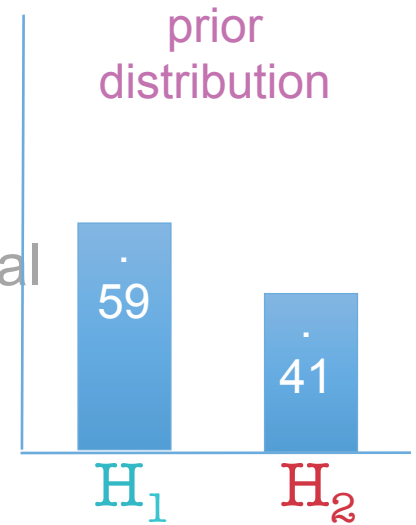
Easier to see in a single update rather than sequential
 Prob. of this example having label Democrat,
 given the values Yes, Yes and No
 on the features NN, TC and LG



$$\text{posterior} \quad P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{
 \begin{array}{c}
 \text{likelihood} \quad P(Y_{\text{NN}} | \text{Dem}) \quad P(Y_{\text{TC}} | \text{Dem}) \quad P(N_{\text{LG}} | \text{Dem}) \quad \text{prior} \quad P(\text{Dem})
 \end{array}
 }{
 P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})
 }$$

We assumed each feature is independent

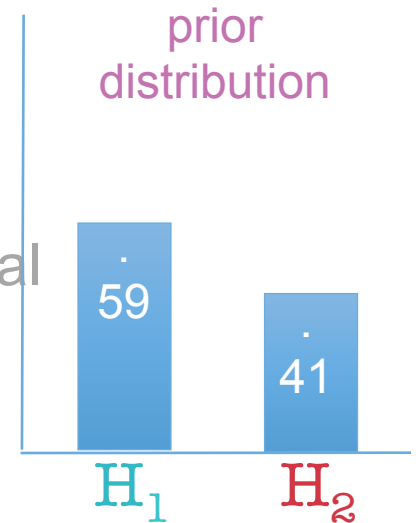
Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG



$$\text{posterior} \quad \text{likelihood} \quad \text{prior}$$
$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{0.949 * 0.169 * 0.898 * 0.59}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

We assumed each feature is independent

Easier to see in a single update rather than sequential
 Prob. of this example having label Democrat,
 given the values Yes, Yes and No
 on the features NN, TC and LG



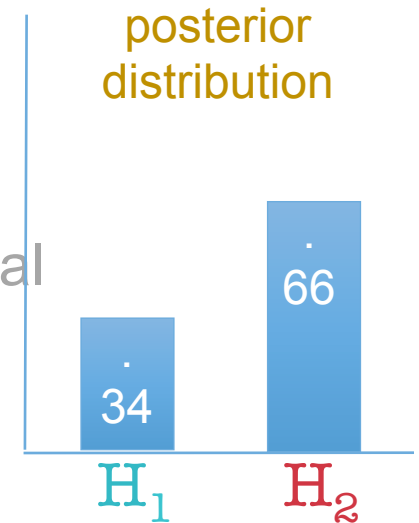
$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{\text{likelihood} \times \text{prior}}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

$$= \frac{0.949 * 0.169 * 0.898 * 0.59}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

$$P(\text{Rep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = \frac{0.829 * 0.854 * 0.561 * 0.41}{P(Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})}$$

We assumed each feature is independent

Easier to see in a single update rather than sequential
Prob. of this example having label Democrat,
given the values Yes, Yes and No
on the features NN, TC and LG



$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = 0.34$$

$$P(\text{Rep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}}) = 0.66 \quad \leftarrow \text{predict!}$$

What about multiple classes?



Straightforward!

Update each hypothesis,
given the values Yes, Yes and No
on the features NN, TC and LG

$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$

$$P(\text{Rep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$

$$P(\text{Indep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$

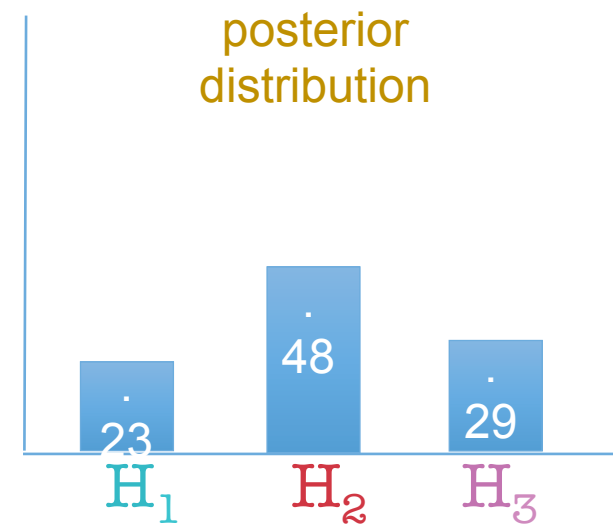
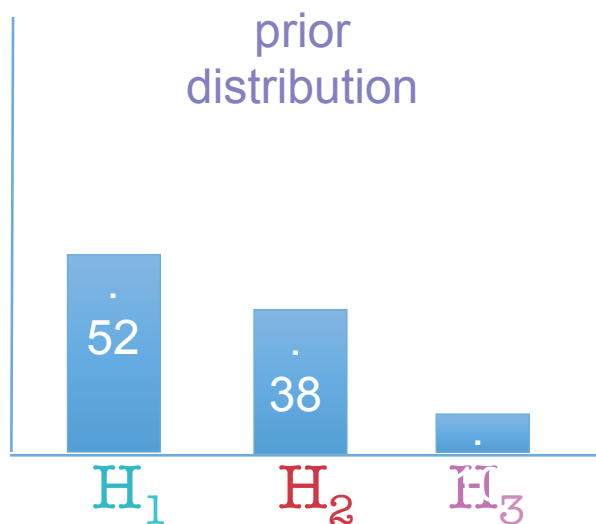
Straightforward!

Update each hypothesis,
given the values **Yes**, **Yes** and **No**
on the features **NN**, **TC** and **LG**

$$P(\text{Dem} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$

$$P(\text{Rep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$

$$P(\text{Indep} | Y_{\text{NN}}, Y_{\text{TC}}, N_{\text{LG}})$$



Naïve Bayes

Bernoulli : models the fraction of documents of class C that contain the word 'w' (ignores number of occurrences)

Vs.

Multinomial: models the fraction of *positions* in documents of class C that contain the word 'w' (keeps track of number of occurrences)

But why does Naïve Bayes work so well-
(Considering that it is Naïve) ?

NB chooses among possible classes to find the class with the highest associated probability.

Naiveté doesn't hurt, because correctness is based on classification, not exact predictions

Advantages of Naïve Bayes:

- Simple & Fast. Just doing a bunch of counts!
- Will converge quickly. Requires less training data
- Can handle sparse matrices
- Can handle multiple classes well

How about numeric features?



5 17 1/4
3 99 28736.123

Naïve Bayes: The Gaussian Approach

$$p(h_x|c) = \frac{1}{\sqrt{2\pi} \sigma_{h,c}^2} \exp - \frac{1}{2} \left(\frac{(h_x - \mu_{h,c})^2}{\sigma_{h,c}^2} \right)$$

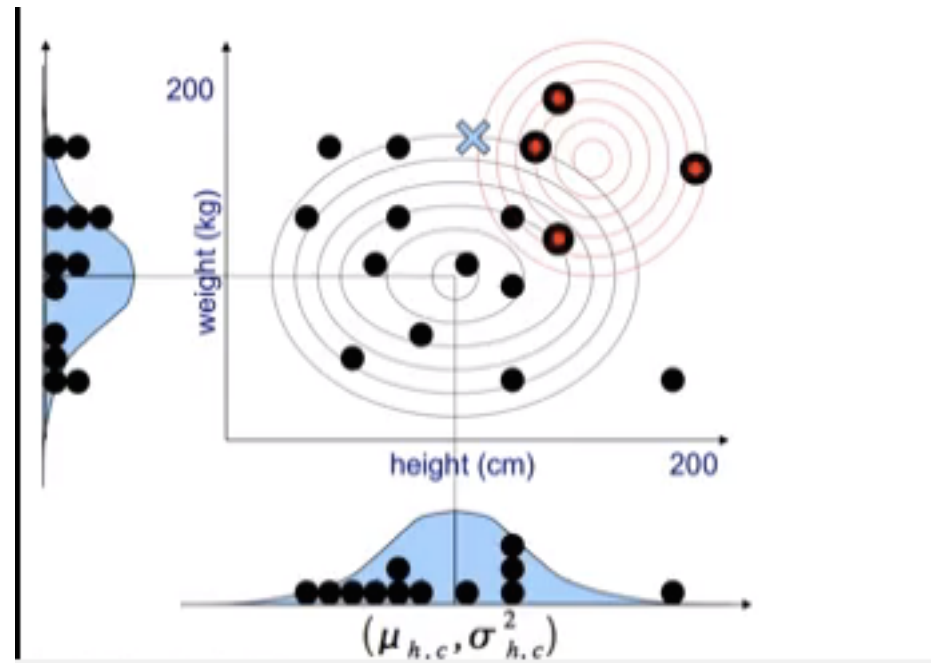
$$p(w_x|c) = \frac{1}{\sqrt{2\pi} \sigma_{w,c}^2} \exp - \frac{1}{2} \left(\frac{(w_x - \mu_{w,c})^2}{\sigma_{w,c}^2} \right)$$

$$p(h_x|a) = \frac{1}{\sqrt{2\pi} \sigma_{h,a}^2} \exp - \frac{1}{2} \left(\frac{(h_x - \mu_{h,a})^2}{\sigma_{h,a}^2} \right)$$

$$p(w_x|a) = \frac{1}{\sqrt{2\pi} \sigma_{w,a}^2} \exp - \frac{1}{2} \left(\frac{(w_x - \mu_{w,a})^2}{\sigma_{w,a}^2} \right)$$

$$P(x|a) = p(h_x|a) p(w_x|a)$$

$$P(a|x) \sim P(x|a) * P(a)$$



Flavors of Bayes in `sklearn`:

Numeric Features:

Gaussian Naïve Bayes

Features that are 0 or 1 (and both matter):

Bernoulli Naïve Bayes

Features that are count-like (and only non-zero matters):

Multinomial Naïve Bayes

Which did we do?