

10-701: Introduction to Machine Learning

Lecture 6 - Naïve Bayes

Henry Chai

2/5/24

Front Matter

- Announcements:

- ~~Nothing!~~ Sign up for AWS credits via Piazza

- Recommended Readings:

- Murphy, Section 3.5

Recall: Coin Flipping

- A Bernoulli random variable takes value **1** (or heads) with probability ϕ and value **0** (or tails) with probability $1 - \phi$
- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

- Assume a Beta prior over the parameter ϕ , which has pdf

$$f(\phi|\alpha, \beta) = \frac{\phi^{\alpha-1}(1 - \phi)^{\beta-1}}{B(\alpha, \beta)}$$

where $B(\alpha, \beta) = \int_0^1 \phi^{\alpha-1}(1 - \phi)^{\beta-1} d\phi$ is a normalizing constant to ensure the distribution integrates to **1**

Example: Beta-Binomial Conjugacy

$$\underbrace{f(\phi|x, \alpha, \beta)}_{\text{Beta posterior}} = \frac{\overset{\text{Bernoulli}}{p(x|\phi)} \overset{\text{Beta prior}}{f(\phi|\alpha, \beta)}}{\underbrace{p(x|\alpha, \beta)}_{\text{evidence}}}$$

$$\begin{aligned} \varphi(x|\alpha, \beta) &= \int_0^1 p(x|\phi) f(\phi|\alpha, \beta) d\phi \\ &= \int_0^1 \phi^x (1-\phi)^{1-x} \frac{\phi^{\alpha-1} (1-\phi)^{\beta-1}}{B(\alpha, \beta)} d\phi \\ &= \frac{1}{B(\alpha, \beta)} \int_0^1 \phi^{\overbrace{(x+\alpha-1)}^{\beta-x}} (1-\phi)^{\overbrace{(1-x+\beta-1)}^{\beta-x}} d\phi \\ &= \frac{B(\alpha+x, \beta-x+1)}{B(\alpha, \beta)} \end{aligned}$$

Example: Beta-Binomial Conjugacy

$$f(\phi|x, \alpha, \beta) = \frac{p(x|\phi)f(\phi|\alpha, \beta)}{p(x|\alpha, \beta)} = \frac{p(x|\phi)f(\phi|\alpha, \beta)}{\left(\frac{B(\alpha+x, \beta-x+1)}{B(\alpha, \beta)}\right)}$$

$$= \frac{\phi^x (1-\phi)^{1-x} \left(\frac{\phi^{\alpha-1} (1-\phi)^{\beta-1}}{B(\alpha, \beta)} \right)}{B(\alpha+x, \beta-x+1)}$$

$$= \frac{\phi^{x+\alpha-1} (1-\phi)^{\beta-x}}{B(\alpha+x, \beta-x+1)}$$

$$= \text{Beta}(\phi | \alpha+x, \beta-x+1)$$

$$= \text{Beta}(\phi | \alpha+x, \beta+(1-x))$$

Beta-Binomial MAP

- Given N iid samples $\{x^{(1)}, \dots, x^{(N)}\}$, the log-posterior is

$$\begin{aligned}\ell(\phi) &= \log f(\phi | \alpha + x^{(1)} + x^{(2)} + \dots + x^{(N)}, \\ &\quad (\beta + (1 - x^{(1)}) + (1 - x^{(2)}) + \dots + (1 - x^{(N)}))) \\ &= \log f(\phi | \alpha + N_1, \beta + N_0)\end{aligned}$$

where N_i is the number of i 's observed in the samples

$$\begin{aligned}&= \log \frac{\phi^{\alpha + N_1 - 1} (1 - \phi)^{\beta + N_0 - 1}}{\cancel{B(\alpha, \beta)} B(\alpha + N_1, \beta + N_0)} \\ &= (\alpha + N_1 - 1) \log \phi + (\beta + N_0 - 1) \log(1 - \phi) - \log \cancel{B(\alpha, \beta)} \\ &\quad B(\alpha + N_1, \beta + N_0)\end{aligned}$$

Beta-Binomial MAP

- Given N iid samples $\{x^{(1)}, \dots, x^{(N)}\}$, the partial derivative of the log-posterior is

$$\frac{\partial \ell}{\partial \phi} = \frac{(\alpha + N_1 - 1)}{\phi} - \frac{(\beta + N_0 - 1)}{1 - \phi}$$
$$\vdots$$

$$\rightarrow \hat{\phi}_{MAP} = \frac{(N_1 + \alpha - 1)}{(N_0 + \beta - 1) + (N_1 + \alpha - 1)}$$

- $\alpha - 1$ is a “pseudocount” of the number of **1**’s you’ve “observed”
- $\beta - 1$ is a “pseudocount” of the number of **0**’s you’ve “observed”

Coin Flipping MAP: Example

- Suppose \mathcal{D} consists of ten 1's or heads ($N_1 = 10$) and two 0's or tails ($N_0 = 2$):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with $\alpha = 2$ and $\beta = 5$, then

$$\phi_{MAP} = \frac{(z - 1 + 10)}{(2 - 1 + 10) + (5 - 1 + 2)} = \frac{11}{17} \approx \frac{1}{2}$$

Coin Flipping MAP: Example

- Suppose \mathcal{D} consists of ten 1's or heads ($N_1 = 10$) and two 0's or tails ($N_0 = 2$):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with $\alpha = 101$ and $\beta = 101$, then

$$\phi_{MAP} = \frac{(101 - 1 + 10)}{(110) + (101 - 1 + 2)} = \frac{110}{212} \approx \frac{1}{2}$$

Coin Flipping MAP: Example

- Suppose \mathcal{D} consists of ten 1's or heads ($N_1 = 10$) and two 0's or tails ($N_0 = 2$):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with $\alpha = 1$ and $\beta = 1$, then

$$\begin{aligned}\phi_{MAP} &= \frac{(1 - 1 + 10)}{(1 - 1 + 10) + (1 - 1 + 2)} = \frac{10}{12} \\ &= \phi_{MLE}\end{aligned}$$

Text Data

- <https://www.nytimes.com/2024/01/30/us/politics/taylor-swift-travis-kelce-trump.html>
- <https://www.breitbart.com/entertainment/2024/01/30/far-left-pro-democrat-facebook-pages-go-all-in-on-taylor-swift-nfl-takeover/>
- https://www.espn.com/nfl/story/_/id/39395830/travis-kelce-taylor-swift-afc-championship/
- <https://www.theonion.com/disillusioned-journalist-begrudgingly-adds-taylor-swift-1850843119>

The New York Times

TaylorSwiftkelce3

FAR LEFT, PRO-DEMOCRAT FACEBOOK PAGES GO ALL IN ON TAYLOR SWIFT NFL TAKEOVER

2 0 2 4 | LIVE Updates | Nevada's Primary and Caucus

Taylor Swift

ESPN

The full girlfriend Kansas C

Share

the ONION

HOME LATEST NEWS LOCAL POLITICS ENTERTAINMENT SPORTS OPINION

Travis Kelce celebrates with Taylor Swift

LOCAL

Disillusioned Journalist Begrudgingly Adds Taylor Swift Reference To Article About Libya Flood

Published September 15, 2023

Twitter Facebook Reddit Email Link

ESPN staff
Jan 28, 2024, 06:33 PM ET

Share Like

The Kansas City Chiefs are headed to the Super Bowl with Taylor Swift in Travis Kelce's family box --

Against the Baltimore Ravens, Kansas City opened the scoring through Kelce, with Patrick Mahomes finding his tight end for a 19-yard score 7 minutes into the first quarter. The Kelce family box celebrated accordingly.

Taylor Swift's vitriol and c

My favorite part about the Super Bowl is how a bunch of people that were supposed to stop watching football over hating will be complaining about Taylor Swift.

Images/Facebook

2:17

r 98%
e right
er Bowl-
enter of

Text Data

- <https://www.nytimes.com/2024/01/30/us/politics/taylor-swift-travis-kelce-trump.html>
- <https://www.breitbart.com/entertainment/2024/01/30/far-left-pro-democrat-facebook-pages-go-all-in-on-taylor-swift-nfl-takeover/>
- https://www.espn.com/nfl/story/_/id/39395830/travis-kelce-taylor-swift-afc-championship/
- <https://www.theonion.com/disillusioned-journalist-begrudgingly-adds-taylor-swift-1850843119>

The New York Times

TaylorSwiftkelce3

FAR LEFT, PRO-DEMOCRAT FACEBOOK PAGES GO ALL IN ON TAYLOR SWIFT NFL TAKEOVER

2024 | LIVE Updates | Nevada's Primary and Caucus

Taylor Swift

ESPN

The full girlfriend Kansas C

Share

the ONION

HOME LATEST NEWS LOCAL POLITICS ENTERTAINMENT SPORTS OPINION

LOCAL

Professor ~~Journalist~~ Begrudgingly Adds Taylor Swift Reference To ~~Article~~ **Lecture**

Published September 15, 2023

Naïve Bayes

Twitter Facebook Reddit Email Link

ESPN staff

Jan 28, 2024, 06:33 PM ET

Share Like

The Kansas City Chiefs are headed to Taylor Swift in Travis Kelce's family box -- v

Against the Baltimore Ravens, Kansas opened the scoring through Kelce, with Patrick Mahomes finding his tight end a 19-yard score 7 minutes into the first quarter. The Kelce family box celebrated accordingly.

Taylor Swift vitriol and c

Images/Facebook

2:17

r 98% e right er Bowl- enter of

Text Data



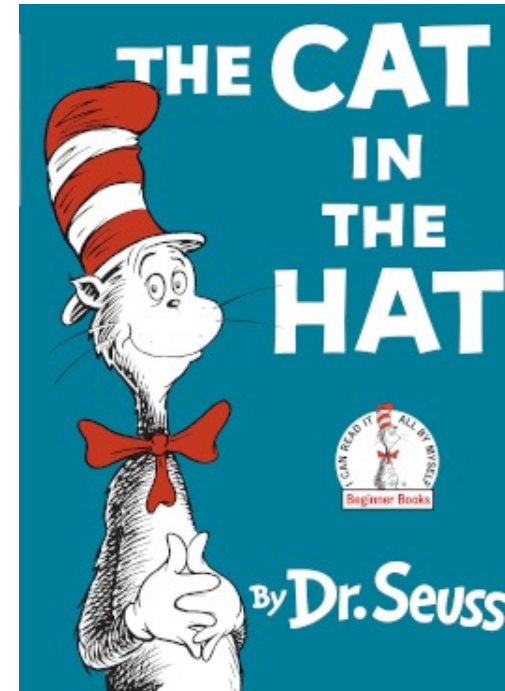
Bag-of-Words Model

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	y (Dr. Seuss)
------------------	------------------	------------------	-------------------	------------------	------------------	--------------------

Bag-of-Words Model

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	y (Dr. Seuss)
1	1	0	0	0	0	1

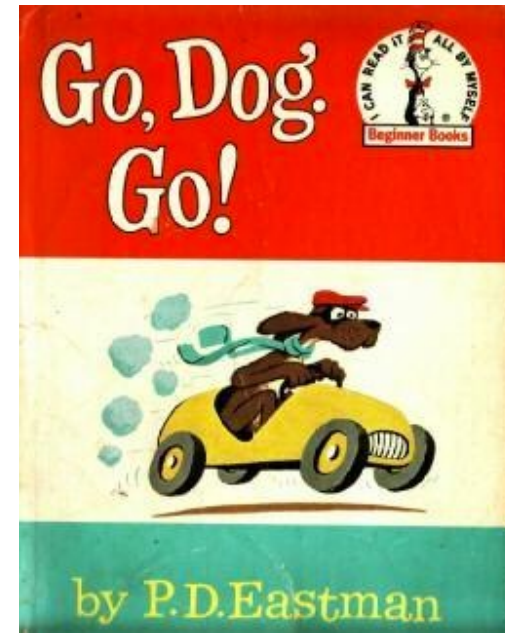
The **Cat** in the **Hat**
(by Dr. Seuss)



Bag-of-Words Model

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	y (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0

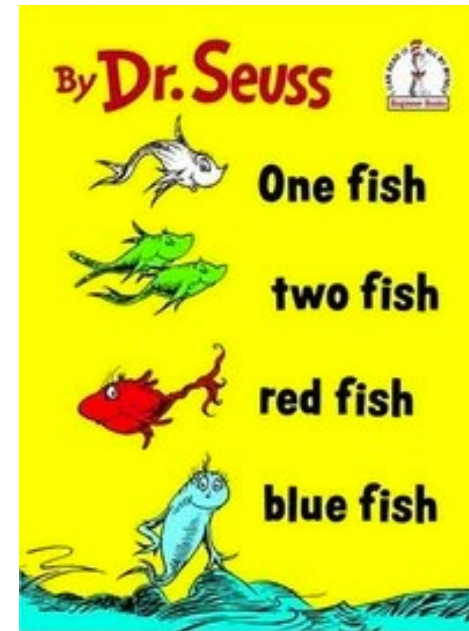
Go, **Dog**. Go!
(by P. D. Eastman)



Bag-of-Words Model

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	y (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1

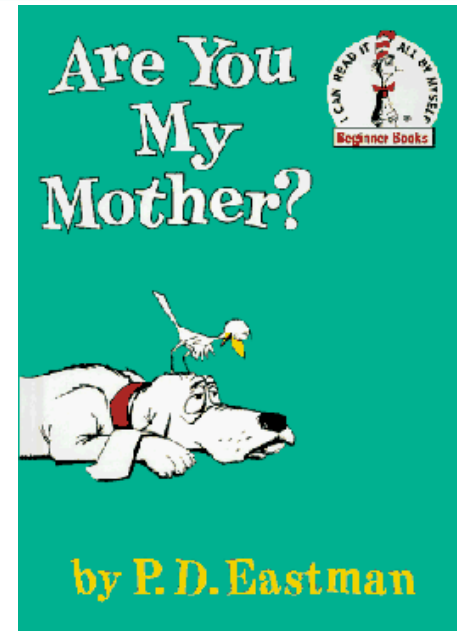
One Fish, Two Fish,
Red Fish, Blue Fish
(by Dr. Seuss)



Bag-of-Words Model

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	y (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1
0	0	0	0	1	0	0

Are You My **Mother**?
(by P. D. Eastman)



Building a Probabilistic Classifier

- Define a decision rule
 - Given a test data point \mathbf{x}' , predict its label \hat{y} using the posterior distribution $P(Y = y|X = \mathbf{x}')$
 - Common choice: $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y|X = \mathbf{x}')$
- Model the posterior distribution
 - Option 1 - Model $P(Y|X)$ directly as some function of X (Wednesday)
 - Option 2 - Use Bayes' rule (today!):

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto P(X|Y) P(Y)$$

How hard is modelling $P(X|Y)$?

- Define a decision rule
 - Given a test data point \mathbf{x}' , predict its label \hat{y} using the posterior distribution $P(Y = y|X = \mathbf{x}')$
 - Common choice: $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y|X = \mathbf{x}')$
- Model the posterior distribution
 - Option 1 - Model $P(Y|X)$ directly as some function of X (later)
 - Option 2 - Use Bayes' rule (today!):

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto P(X|Y) P(Y)$$

$P(x_1=1 | Y=1)$; I don't need to learn
 $\nearrow P(x_1=1 | Y=0)$ $P(x_1=0 | Y=1)$ or $P(x_1=0 | Y=0)$

x_1 ("hat")	x_2 ("cat")	x_3 ("dog")	x_4 ("fish")	x_5 ("mom")	x_6 ("dad")	$P(X Y = 1)$
0	0	0	0	0	0	θ_1
1	0	0	0	0	0	θ_2
1	1	0	0	0	0	θ_3
1	0	1	0	0	0	θ_4

How hard is
modelling
 $P(X|Y)$?

Using the
Naive Bayes
assumption?

Naïve Bayes Assumption

- **Assume** features are conditionally independent given the label:

$$P(x|Y) = P(x_1 \cap x_2 \cap \dots \cap x_D | Y) \\ = \prod_{d=1}^D P(x_d | Y)$$

- Pros:

★ — significantly reduces the # of parameters

— reduces overfitting!

- Cons:

— mathematically convenient

— This a terrible assumption!

— \exists a relaxed of this assumption
 \Rightarrow Bayesian networks

General Recipe for Machine Learning

- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for Naïve Bayes

$$\vec{\Theta} = [\theta_{1,0}, \theta_{1,1}, \theta_{2,0}, \theta_{2,1}, \dots]$$

- Define a model and model parameters

- Make the Naïve Bayes assumption

- Assume iid data points

- Parameters: $\theta_{d,y} = P(X_d = 1 | Y = y)$; $\pi = P(Y = 1)$

- Write down an objective function

- Maximize the log-likelihood

- Optimize the objective w.r.t. the model parameters

- Solve in closed-form

Setting the Parameters via MLE

$$\begin{aligned} \ell_D(\pi, \Theta) &= \log \left(\frac{P(D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\})}{|\pi, \Theta|} \right) \\ &= \log \prod_{n=1}^N P(x^{(n)}, y^{(n)} | \pi, \Theta) \\ &= \log \prod_{n=1}^N \left(\underbrace{P(x^{(n)} | y^{(n)}, \Theta)}_{\substack{D \\ \theta_{d,1}, \theta_{d,0}}} P(y^{(n)} | \pi) \right) \\ &= \log \prod_{n=1}^N \left(\prod_{d=1}^D P(x_d^{(n)} | y^{(n)}, \theta_{d,1}, \theta_{d,0}) \right) P(y^{(n)} | \pi) \\ &= \sum_{n=1}^N \sum_{d=1}^D \log \left(P(x_d^{(n)} | y^{(n)}, \theta_{d,1}, \theta_{d,0}) \right) + \log P(y^{(n)} | \pi) \\ &= \sum_{n: y^{(n)}=1} \left(\sum_{d=1}^D \log \left(P(x_d^{(n)} | y^{(n)}=1, \theta_{d,1}) \right) + \log P(y^{(n)}=1 | \pi) \right) \\ &\quad + \sum_{n: y^{(n)}=0} \left(\dots y^{(n)}=0, \theta_{d,0} \right) \end{aligned}$$

Setting the Parameters via MLE

- Binary label
 - $Y \sim \text{Bernoulli}(\pi)$
 - $\hat{\pi} = N_{Y=1} / N$
 - $N = \#$ of data points
 - $N_{Y=1} = \#$ of data points with label 1
- Binary features
 - $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$
 - $\hat{\theta}_{d,y} = N_{Y=y, X_d=1} / N_{Y=y}$
 - $N_{Y=y} = \#$ of data points with label y
 - $N_{Y=y, X_d=1} = \#$ of data points with label y and feature $X_d = 1$

Bernoulli Naïve Bayes

- Binary label
 - $Y \sim \text{Bernoulli}(\pi)$
 - $\hat{\pi} = N_{Y=1} / N$
 - $N = \#$ of data points
 - $N_{Y=1} = \#$ of data points with label 1
- Binary features
 - $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$
 - $\hat{\theta}_{d,y} = N_{Y=y, X_d=1} / N_{Y=y}$
 - $N_{Y=y} = \#$ of data points with label y
 - $N_{Y=y, X_d=1} = \#$ of data points with label y and feature $X_d = 1$

Multiclass Bernoulli Naïve Bayes

- Discrete label (Y can take on one of M possible values)
 - $Y \sim \text{Categorical}(\pi_1, \dots, \pi_M)$
 - $\hat{\pi}_m = N_{Y=m} / N$
 - $N = \#$ of data points
 - $N_{Y=m} = \#$ of data points with label m
- Binary features
 - $X_d | Y = m \sim \text{Bernoulli}(\theta_{d,m})$
 - $\hat{\theta}_{d,m} = N_{Y=m, X_d=1} / N_{Y=m}$
 - $N_{Y=m} = \#$ of data points with label m
 - $N_{Y=m, X_d=1} = \#$ of data points with label m and feature $X_d = 1$

Multinomial Naïve Bayes

- Binary label
 - $Y \sim \text{Bernoulli}(\pi)$
 - $\hat{\pi} = N_{Y=1} / N$
 - $N = \#$ of data points
 - $N_{Y=1} = \#$ of data points with label 1
- Discrete features (X_d can take on one of K possible values)
 - $X_d | Y = y \sim \text{Categorical}(\theta_{d,1,y}, \dots, \theta_{d,K,y})$
 - $\hat{\theta}_{d,k,y} = N_{Y=y, X_d=k} / N_{Y=y}$
 - $N_{Y=y} = \#$ of data points with label y
 - $N_{Y=y, X_d=k} = \#$ of data points with label y and feature $X_d = k$

Gaussian Naïve Bayes

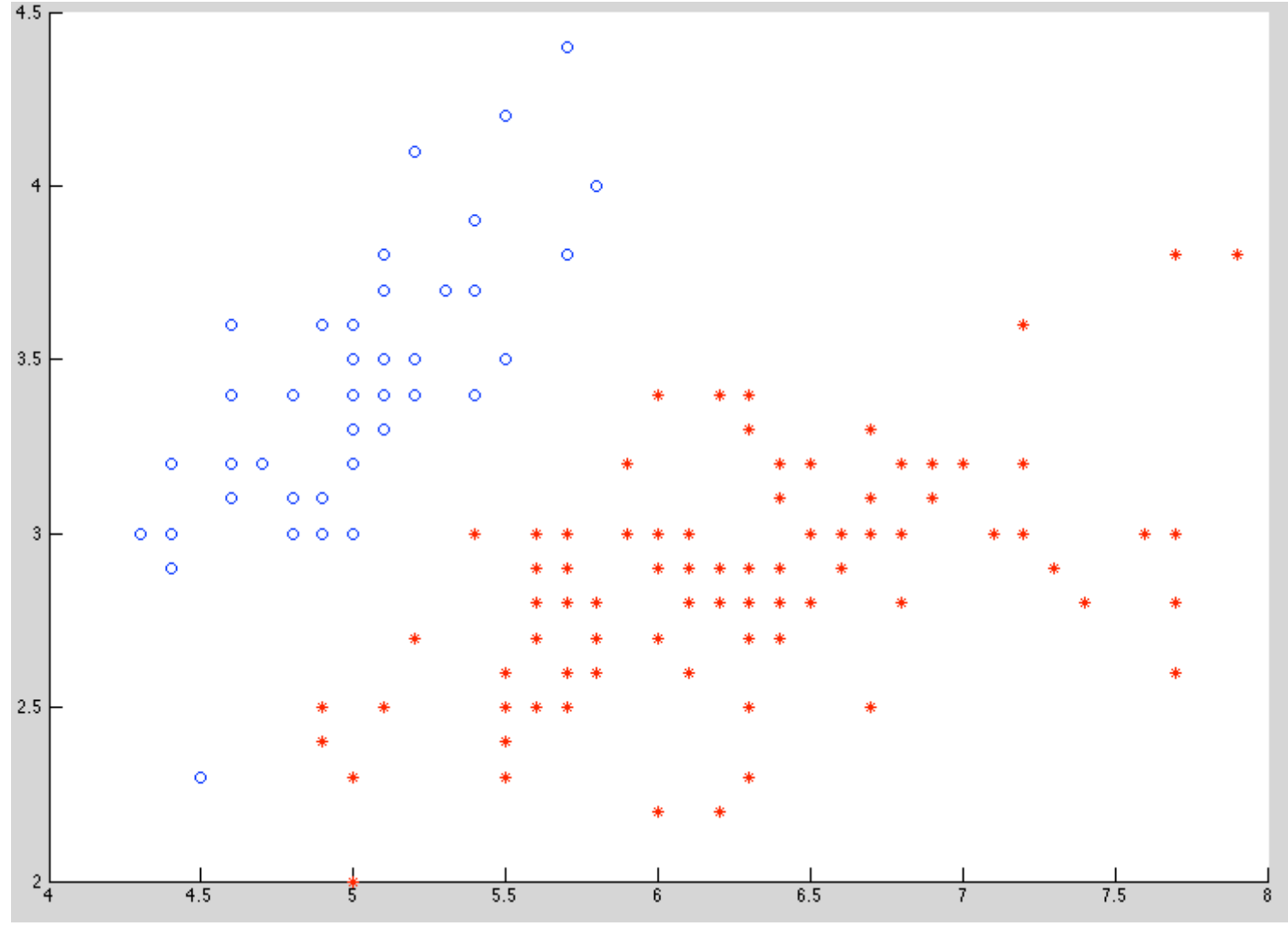
- Binary label
 - $Y \sim \text{Bernoulli}(\pi)$
 - $\hat{\pi} = N_{Y=1} / N$
 - $N = \#$ of data points
 - $N_{Y=1} = \#$ of data points with label 1
- Real-valued features
 - $X_d | Y = y \sim \text{Gaussian}(\mu_{d,y}, \sigma_{d,y}^2)$
 - $\hat{\mu}_{d,y} = \frac{1}{N_{Y=y}} \sum_{n:y^{(n)}=y} x_d^{(n)}$
 - $\hat{\sigma}_{d,y}^2 = \frac{1}{N_{Y=y}} \sum_{n:y^{(n)}=y} (x_d^{(n)} - \hat{\mu}_{d,y})^2$
 - $N_{Y=y} = \#$ of data points with label y

Recall: Fisher Iris Dataset

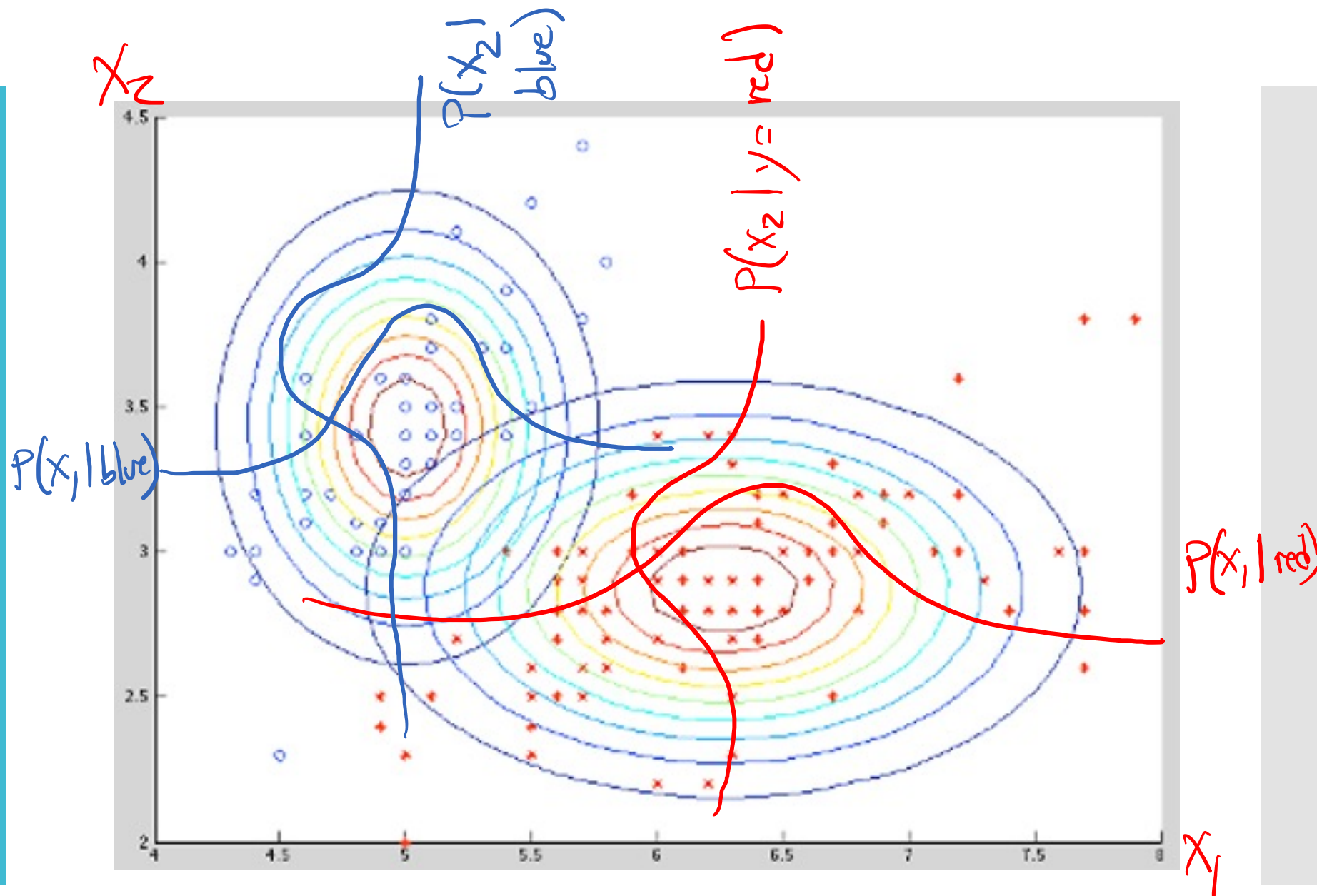
- Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width
0	4.3	3.0
0	4.9	3.6
0	5.3	3.7
1	4.9	2.4
1	5.7	2.8
1	6.3	3.3
1	6.7	3.0

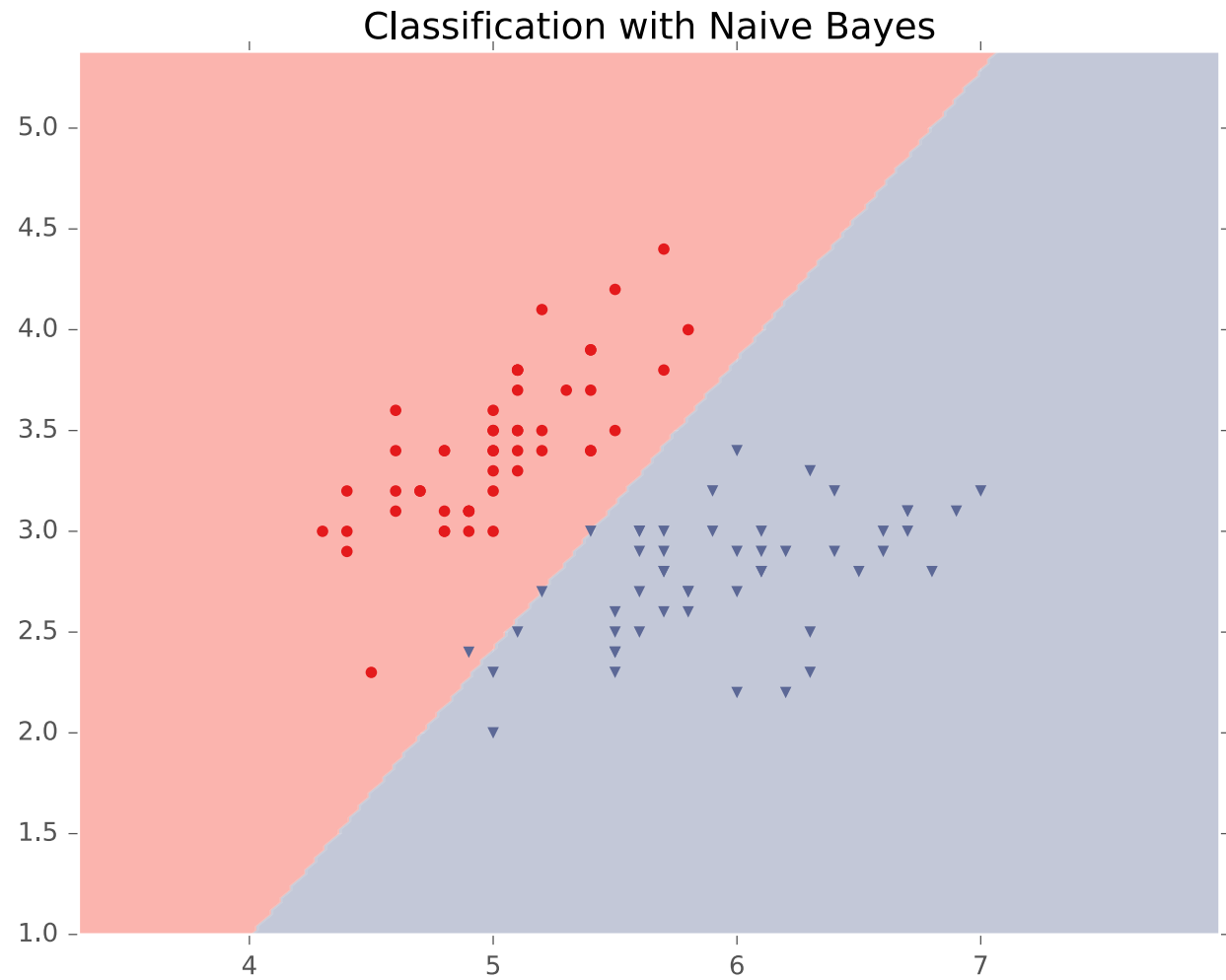
Visualizing Gaussian Naïve Bayes (2 classes)



Visualizing Gaussian Naïve Bayes (2 classes)



Visualizing
Gaussian
Naïve
Bayes
(2 classes,
equal
variances)



Visualizing Gaussian Naïve Bayes (2 classes, learned variances)

