

10-301/601: Introduction to Machine Learning

Lecture 16 - Naïve Bayes

Henry Chai & Matt Gormley

10/25/23

Front Matter

- Announcements:
 - HW5 released 10/9, due 10/27 (Friday) at 11:59 PM
 - HW6 released 10/27 (Friday), due 11/3 at 11:59 PM
 - **You may only use at most 2 late days on HW6**
 - Exam 2 on 11/9
 - All topics between Lecture 8 and Lecture 16 (today's lecture) are in-scope
 - Exam 1 content may be referenced but will not be the primary focus of any question
 - Exam 3 on 12/12 from 5:30 PM to 8:30 PM
 - Sign up for peer tutoring! See [Piazza](#) for more details

Recall: Maximum a Posteriori (MAP) Estimation

- Insight: sometimes we have *prior* information we want to incorporate into parameter estimation
- Idea: use Bayes rule to reason about the *posterior* distribution over the parameters

- MLE finds $\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)$

- MAP finds $\hat{\theta} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{D})$ ← posterior
 $= \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)p(\theta)/p(\mathcal{D})$

$$= \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)p(\theta)$$

likelihood prior

$$= \operatorname{argmax}_{\theta} \underbrace{\log p(\mathcal{D}|\theta) + \log p(\theta)}_{\text{log-posterior}}$$

Recall: Maximum a Posteriori (MAP) Estimation

1. Specify the generative story, i.e., the data generating distribution, including a prior distribution

2. Maximize the log-posterior of $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$

$$\ell_{MAP}(\theta) = \log p(\theta) + \sum_{i=1}^N \log p(x^{(i)} | \theta)$$

3. Solve in closed form: take partial derivatives, set to 0 and solve

Coin Flipping MAP

- 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

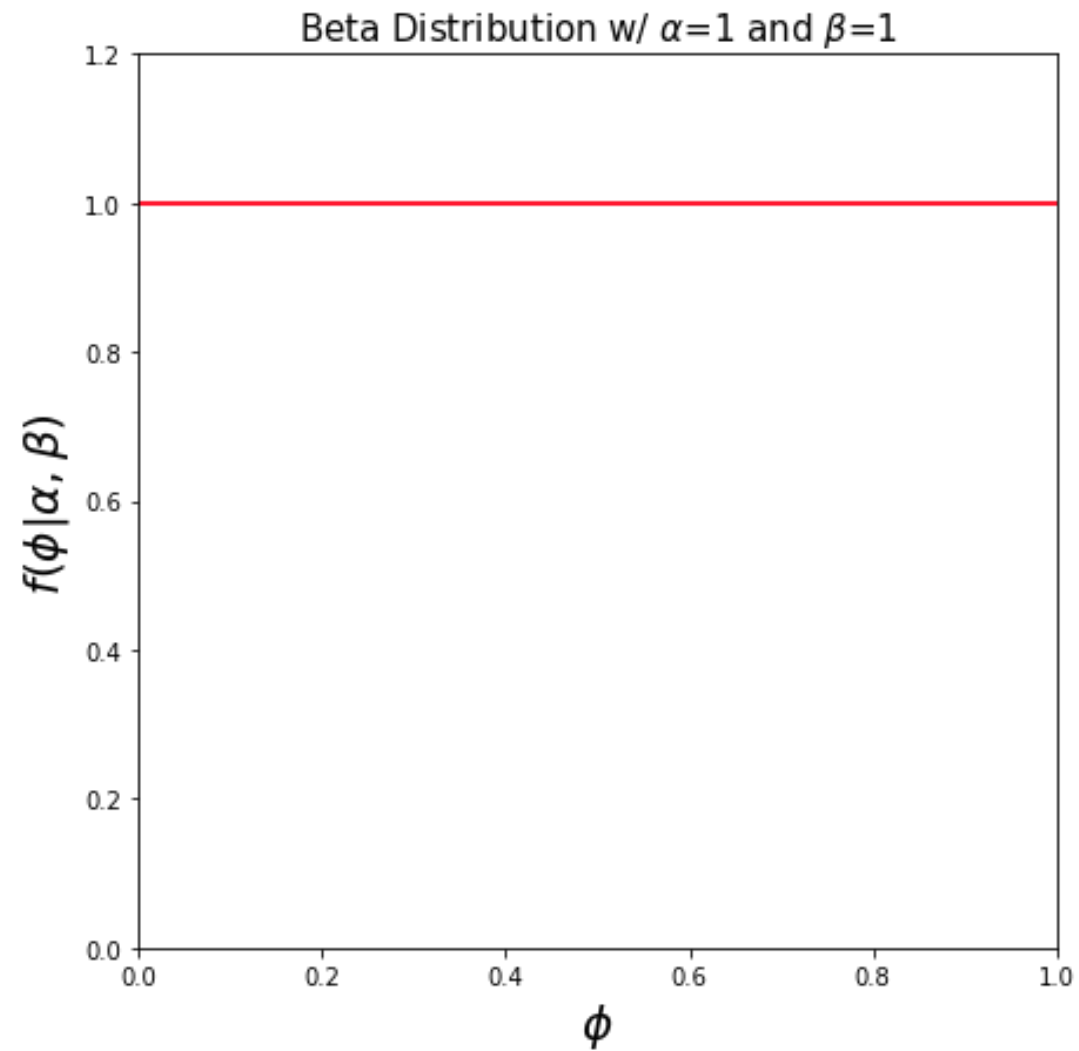
$$\longrightarrow p(x|\phi) = \phi^x (1 - \phi)^{1-x}$$

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

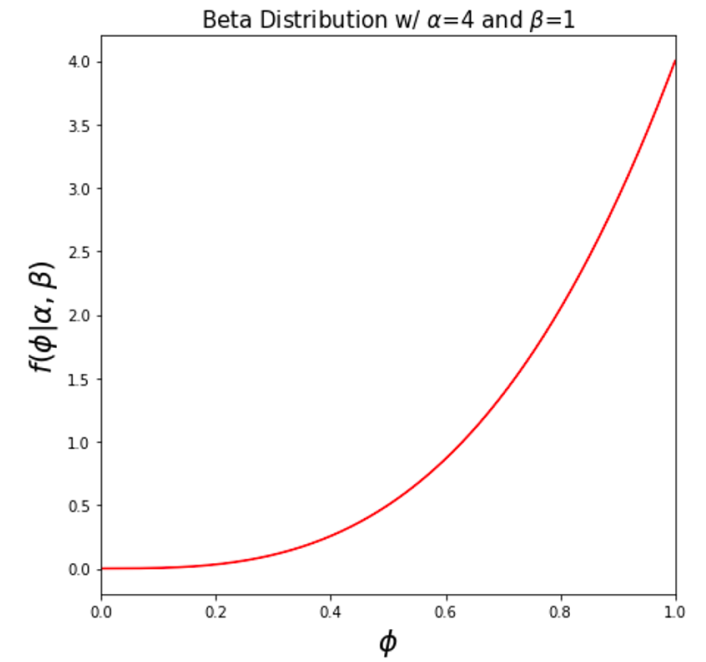
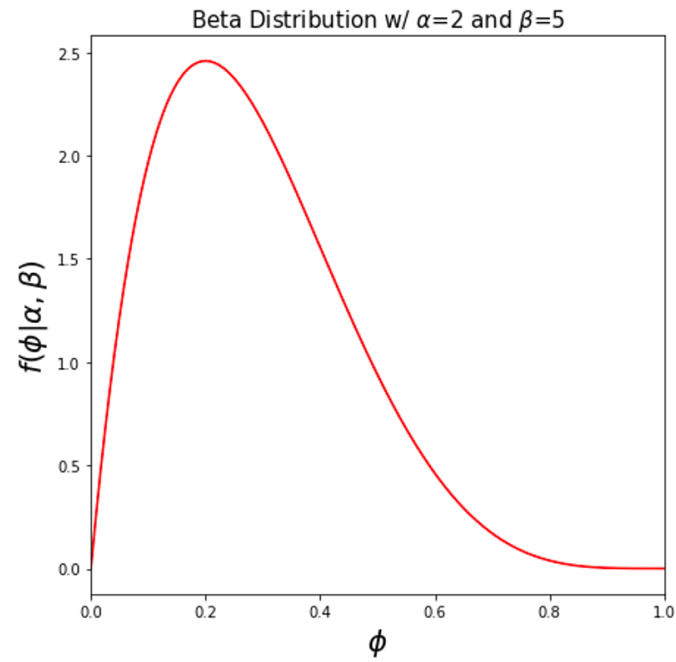
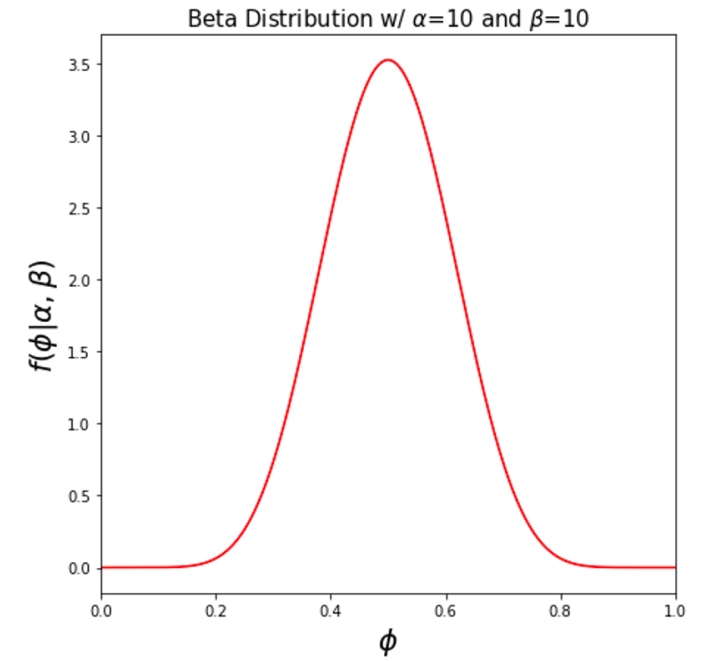
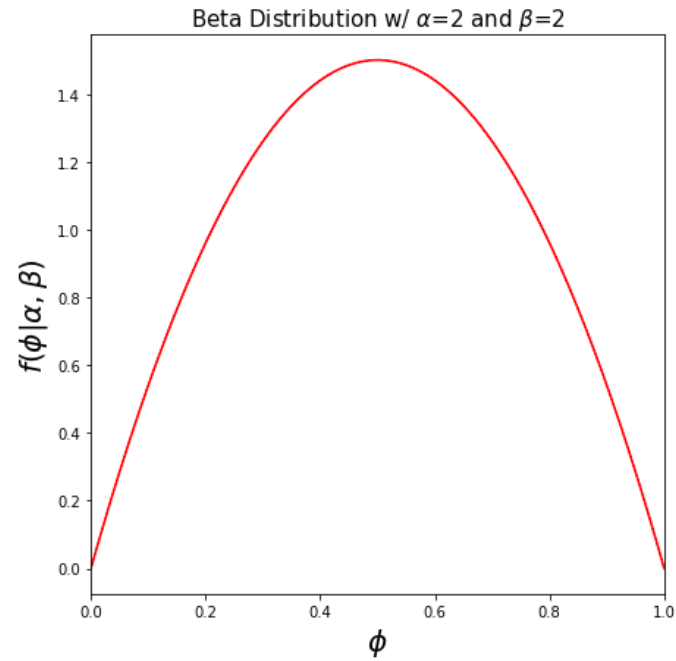
$$f(\phi|\alpha, \beta) = \frac{\phi^{\alpha-1} (1 - \phi)^{\beta-1}}{\underline{B(\alpha, \beta)}} \quad \longleftarrow$$

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**

Beta Distribution



Beta Distribution



Why use this strange looking Beta prior?

The Beta distribution is the *conjugate prior* for the Bernoulli distribution!

- 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

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

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**

Coin Flipping MAP

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

$$l(\phi) = \log(p(\phi | D)) \propto \log(p(\phi)) + \log(p(D | \phi))$$

$$= \log\left(\frac{\phi^{\alpha-1} (1-\phi)^{\beta-1}}{B(\alpha, \beta)}\right) + \underbrace{N_1 \log \phi + N_0 \log(1-\phi)}$$

$$= (\alpha-1) \log \phi + (\beta-1) \log(1-\phi) - \log(B(\alpha, \beta))$$

$$+ N_1 \log \phi + N_0 \log(1-\phi)$$

$$\rightarrow = (\alpha-1 + N_1) \log \phi + (\beta-1 + N_0) \log(1-\phi)$$

$$- \log(B(\alpha, \beta))$$

Coin Flipping MAP

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

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

⋮

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


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

$$\phi_{MAP} = \frac{(2 - 1 + 10)}{(2 - 1 + 10) + (5 - 1 + 2)} = \frac{11}{17} < \frac{10}{12}$$


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

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

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

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

MLE/MAP Learning Objectives

You should be able to...

- Recall probability basics, including but not limited to: discrete and continuous random variables, probability mass functions, probability density functions, events vs. random variables, expectation and variance, joint probability distributions, marginal probabilities, conditional probabilities, independence, conditional independence
- State the principle of maximum likelihood estimation and explain what it tries to accomplish
- State the principle of maximum a posteriori estimation and explain why we use it
- Derive the MLE or MAP parameters of a simple model in closed form

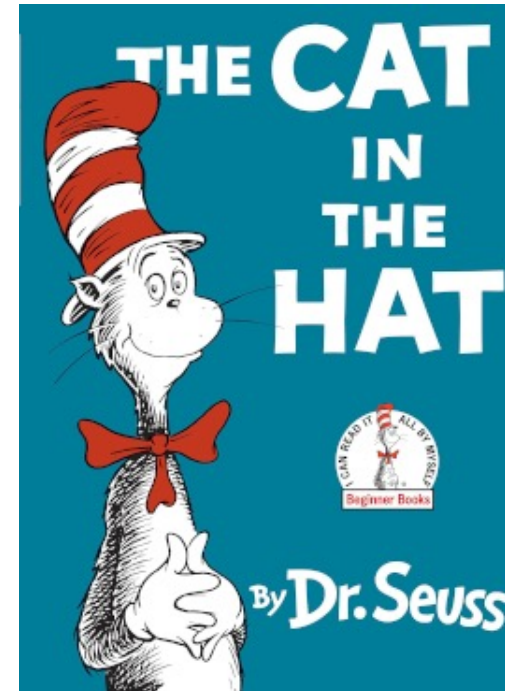
Bag-of-Words Model

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

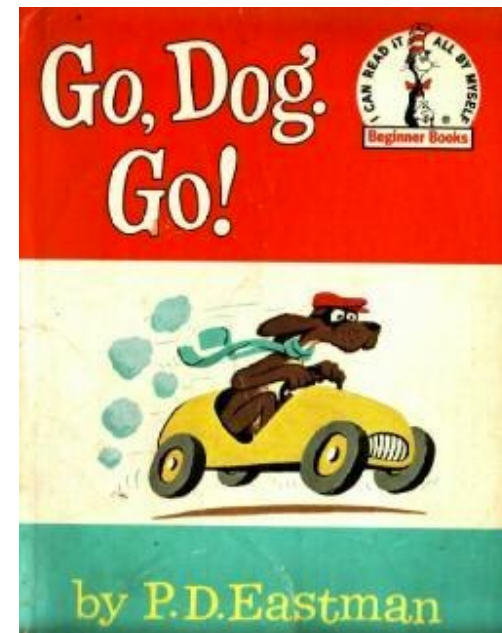
The **Cat** in the **Hat**
(by Dr. Seuss)



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

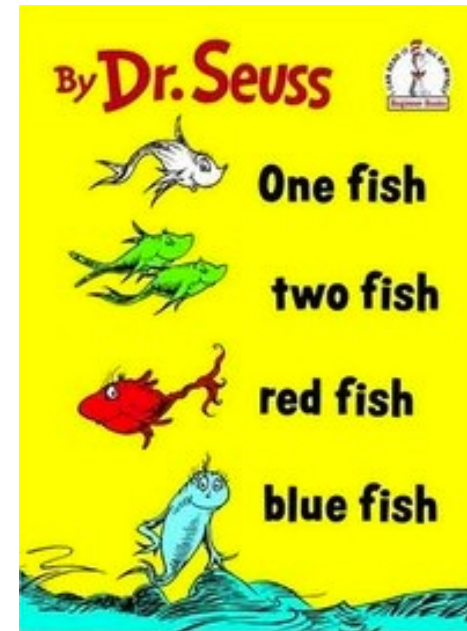
Go, **Dog**. Go!
(by P. D. Eastman)



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

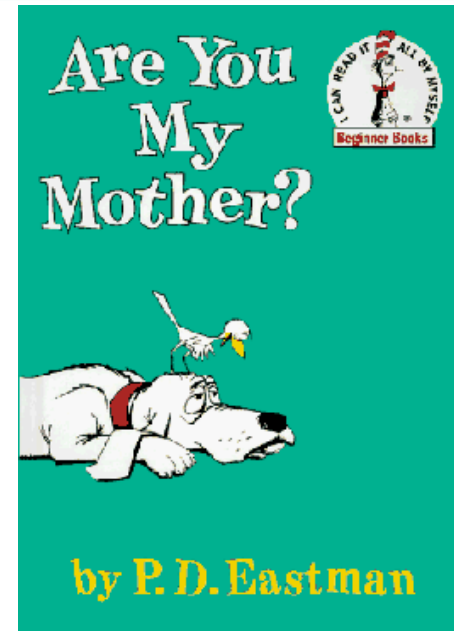
One Fish, Two Fish,
Red Fish, Blue Fish
(by Dr. Seuss)



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 (recall: logistic regression)
 - Option 2 - Use Bayes' rule:

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

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 (recall: logistic regression)
 - Option 2 - Use Bayes' rule:

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

$$\begin{aligned} & \rightarrow P(x_1=1|Y=1) \Rightarrow P(x_1=0|Y=1) = 1 - P(x_1=1|Y=1) \\ & \rightarrow P(x_1=1|Y=0) \end{aligned}$$

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 Naive
Bayes

Naïve Bayes Assumption

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

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

- Pros:

- Significantly reduces # of parameters
- Can help combat overfitting

- Cons:

- This is like a bad assumption
or incorrect

- Can relax this assumption \Rightarrow Bayesian Network

Recipe for Naïve Bayes

- Define a model and model parameters
 - make this Naïve Bayes assumption
 - Assume iid data $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$
 - Parameters: $\Theta_{d,y} = P(X_d = 1 | Y = y)$, $\pi = P(Y = 1)$
- Write down an objective function
 - Maximizing the log-likelihood
- Optimize the objective w.r.t. the model parameters
 - Solve for these in closed-form

Setting the Parameters via MLE

$$\ell_{\mathcal{D}}(\pi, \theta) = \log P(\mathcal{D} = \{x^{(1)}, y^{(1)}, \dots, x^{(N)}, y^{(N)}\} | \pi, \theta)$$

$$= \log \prod_{i=1}^N P(\vec{x}^{(i)}, y^{(i)} | \pi, \theta)$$

$$= \log \prod_{i=1}^N \underbrace{P(\vec{x}^{(i)} | y^{(i)}, \theta)}_{\text{conditional}} P(y^{(i)} | \pi)$$

$$= \log \prod_{i=1}^N \left(\prod_{d=1}^D P(x_d^{(i)} | y^{(i)}, \theta_{d,1}, \theta_{d,0}) \right) P(y^{(i)} | \pi)$$

$$= \sum_{i=1}^N \left(\sum_{d=1}^D \log P(x_d^{(i)} | y^{(i)}, \theta_{d,1}, \theta_{d,0}) \right) + \log P(y^{(i)} | \pi)$$

$$= \sum_{\substack{y^{(i)}=1 \\ y^{(i)}=0}} \left(\sum_{d=1}^D \log P(x_d^{(i)} | \theta_{d,1}) \right) + \sum_{\substack{y^{(i)}=0 \\ y^{(i)}=0}} \left(\sum_{d=1}^D \log P(x_d^{(i)} | \theta_{d,0}) \right) + \sum_{i=1}^N \log P(y^{(i)} | \pi)$$

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 $X_d = 1$

Poll Question 1:
Given this
dataset, what is
the MLE of π ?

Poll Question 2:
Given this
dataset, what is
the MLE of $\theta_{3,1}$?

x_1	x_2	x_3	y
1	0	1	0
0	1	0	1
0	1	1	1
0	0	1	0
1	0	1	0
1	0	1	1

A. 0/6

B. 1/6

C. 2/6

1. D. 3/6

2. E. 4/6

F. 5/6

G. 6/6

H. 7/6 (TOXIC)

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 $X_d = 1$

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)
 - $X_d | Y = y \sim \text{Categorical}(\theta_{d,1,y}, \dots, \theta_{d,K-1,y})$
 - $\hat{\theta}_{d,k,y} = N_{Y=y, X_d=k} / N_{Y=y}$
 - $N_{Y=y} = \#$ 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
 - $X_d | Y = y \sim \text{Gaussian}(\mu_{d,y}, \sigma_{d,y}^2)$
 - • $\hat{\mu}_{d,y} = \frac{1}{N_{Y=y}} \sum_{n:y^{(n)}=y} x_d^{(n)}$
 - • $\hat{\sigma}_{d,y}^2 = \frac{1}{N_{Y=y}} \sum_{n:y^{(n)}=y} \left(x_d^{(n)} - \hat{\mu}_{d,y} \right)^2$
 - $N_{Y=y} = \#$ of data points with label y

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

Visualizing Gaussian Naïve Bayes

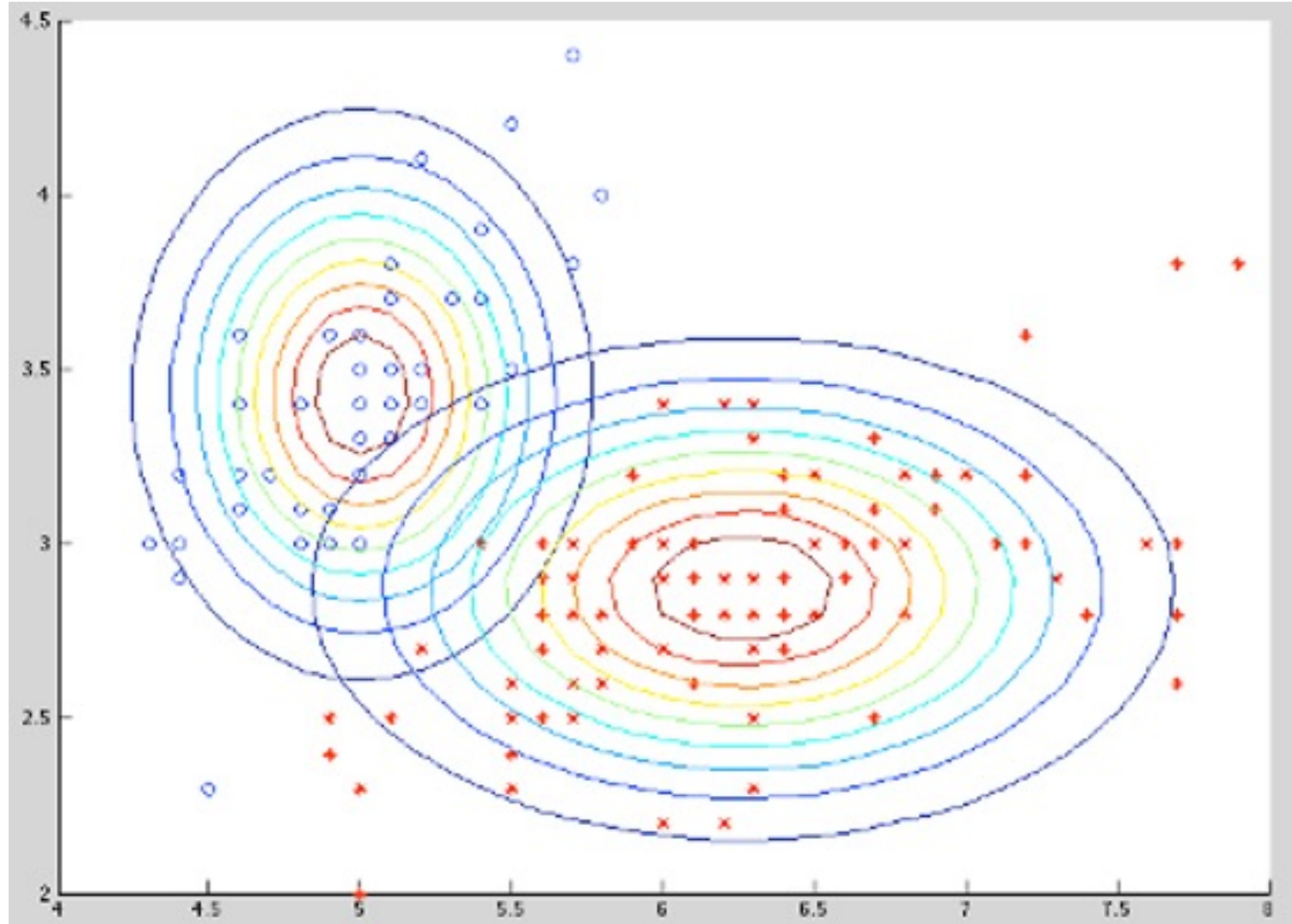
- Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

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

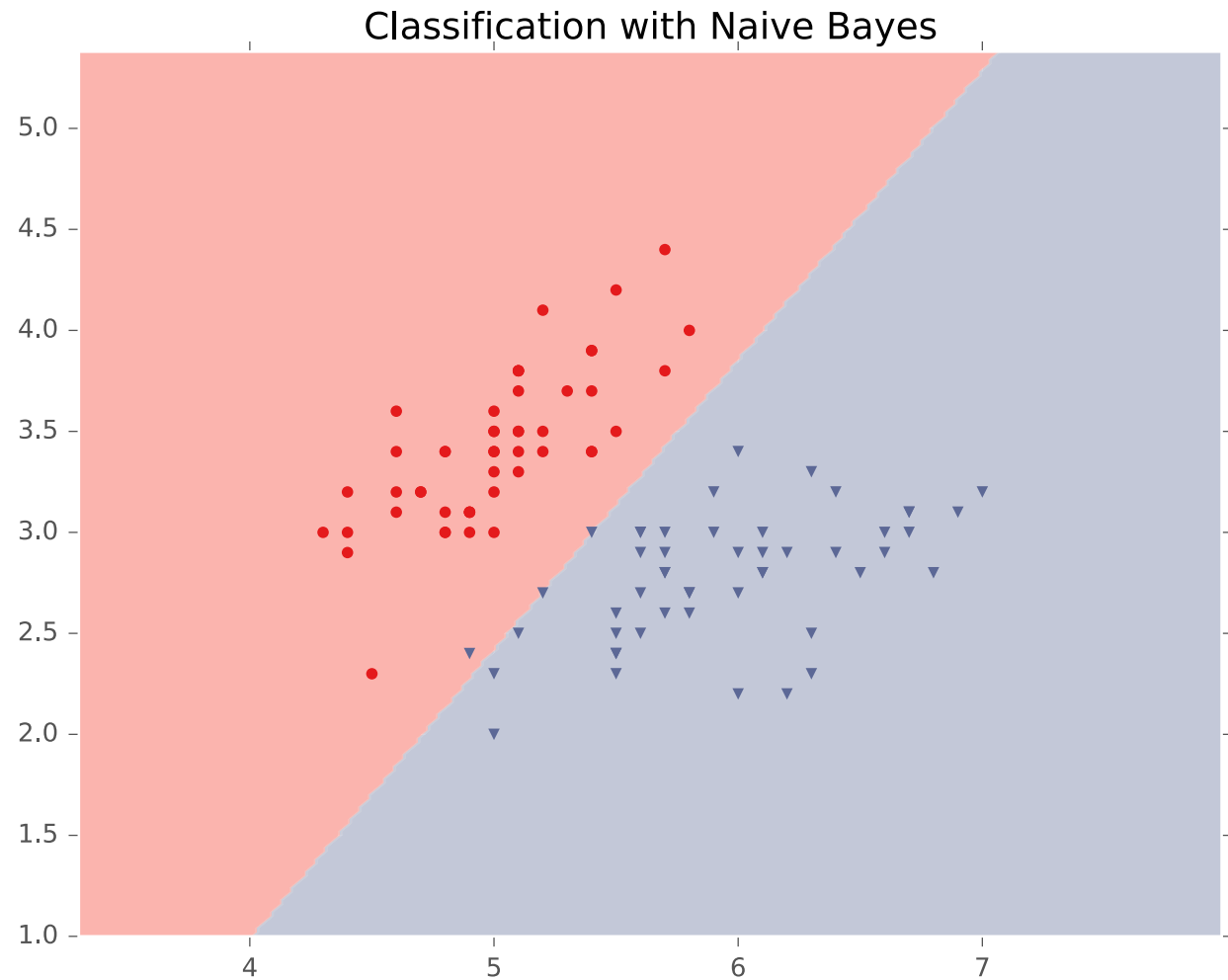
Deleted two of the four features, so that input space is 2D



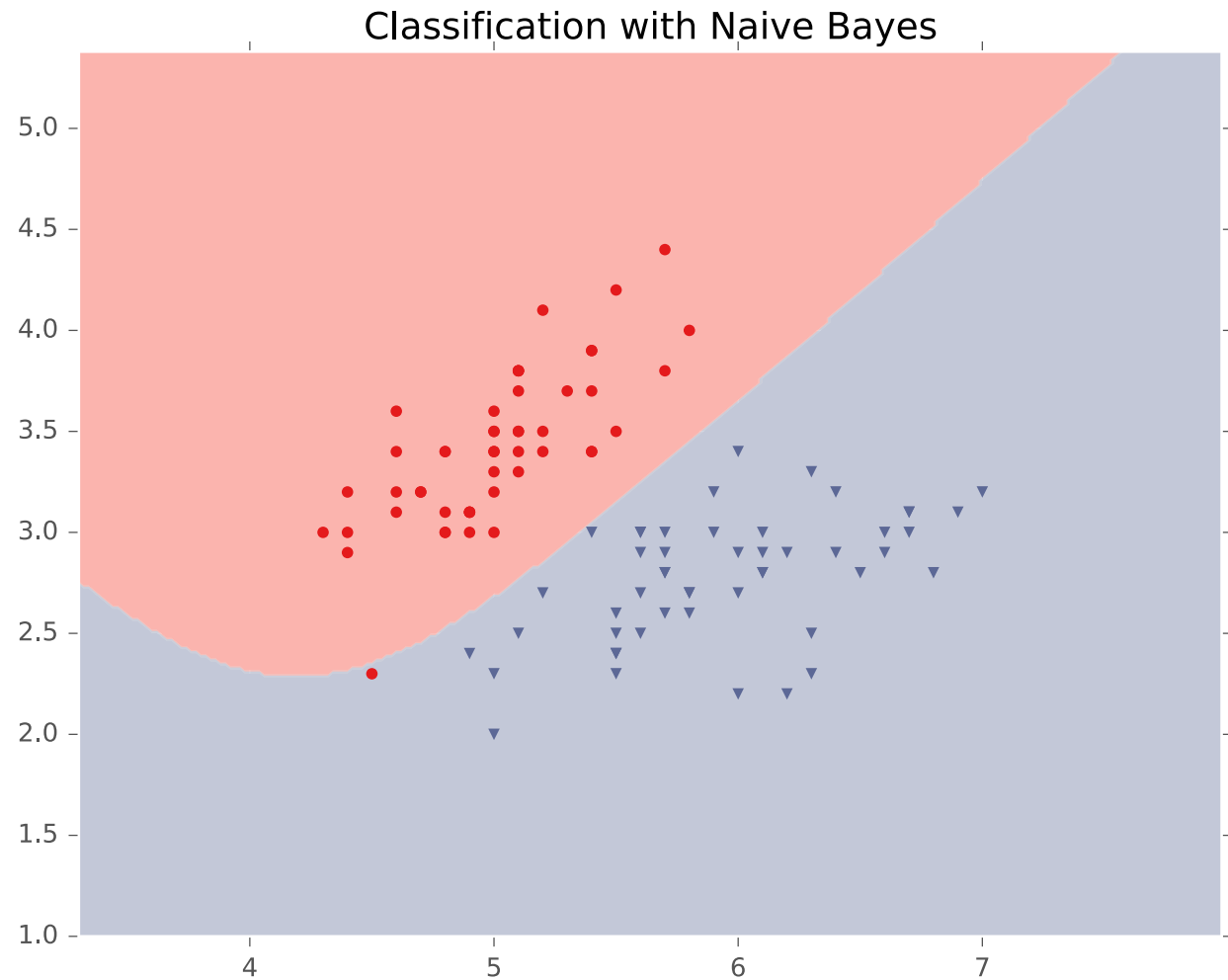
Visualizing Gaussian Naïve Bayes (2 classes)



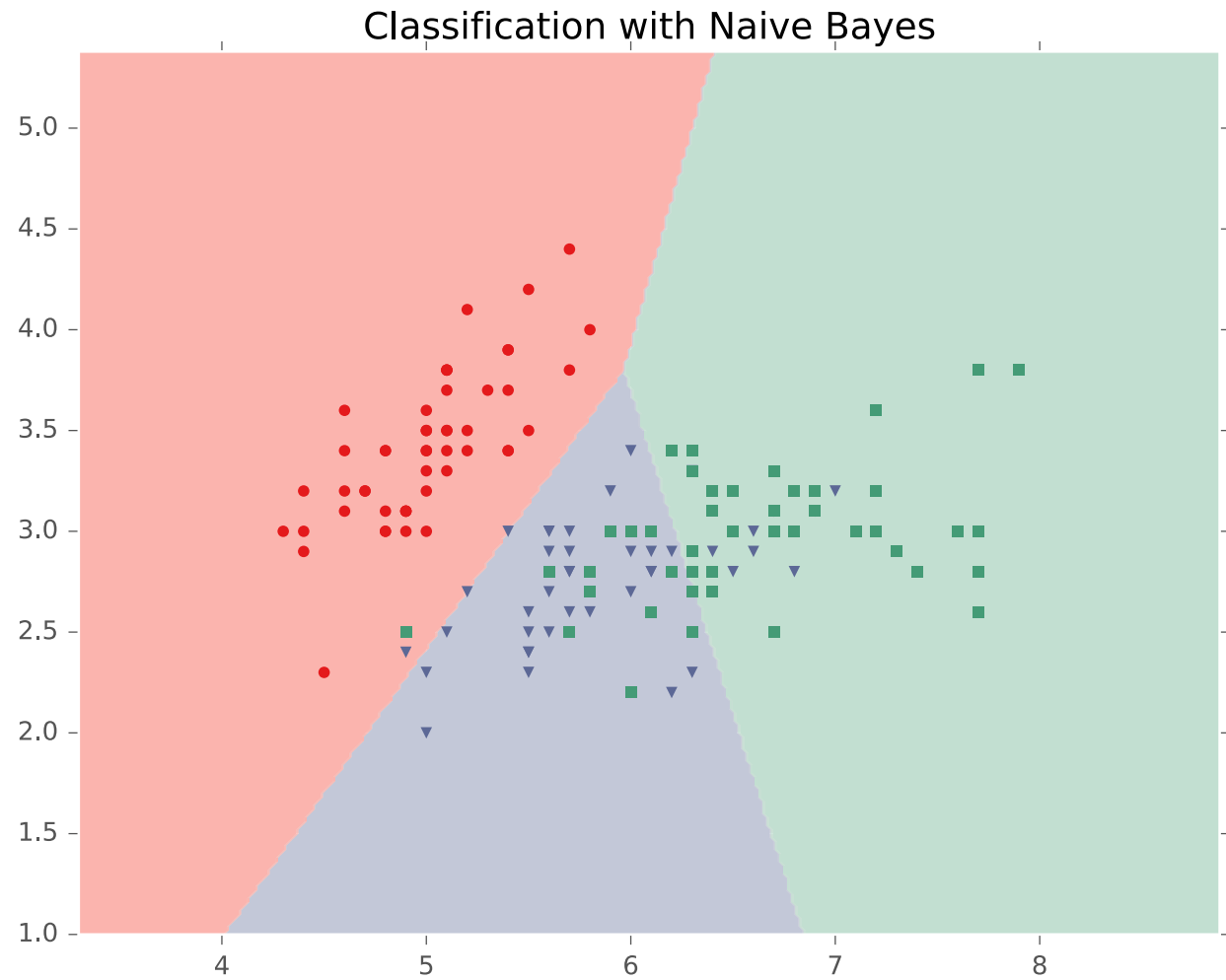
Visualizing Gaussian Naïve Bayes (2 classes, equal variances)



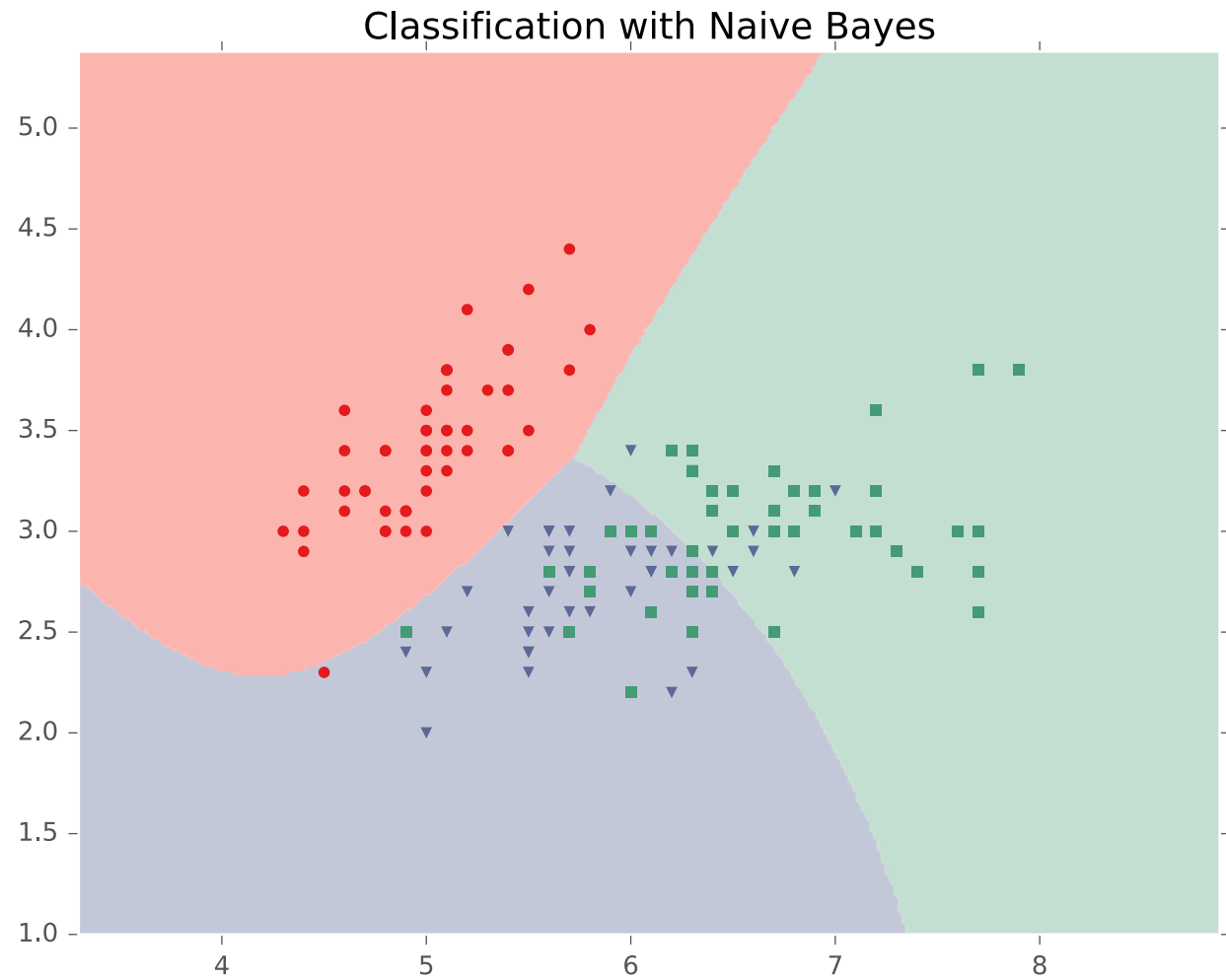
Visualizing Gaussian Naïve Bayes (2 classes, learned variances)



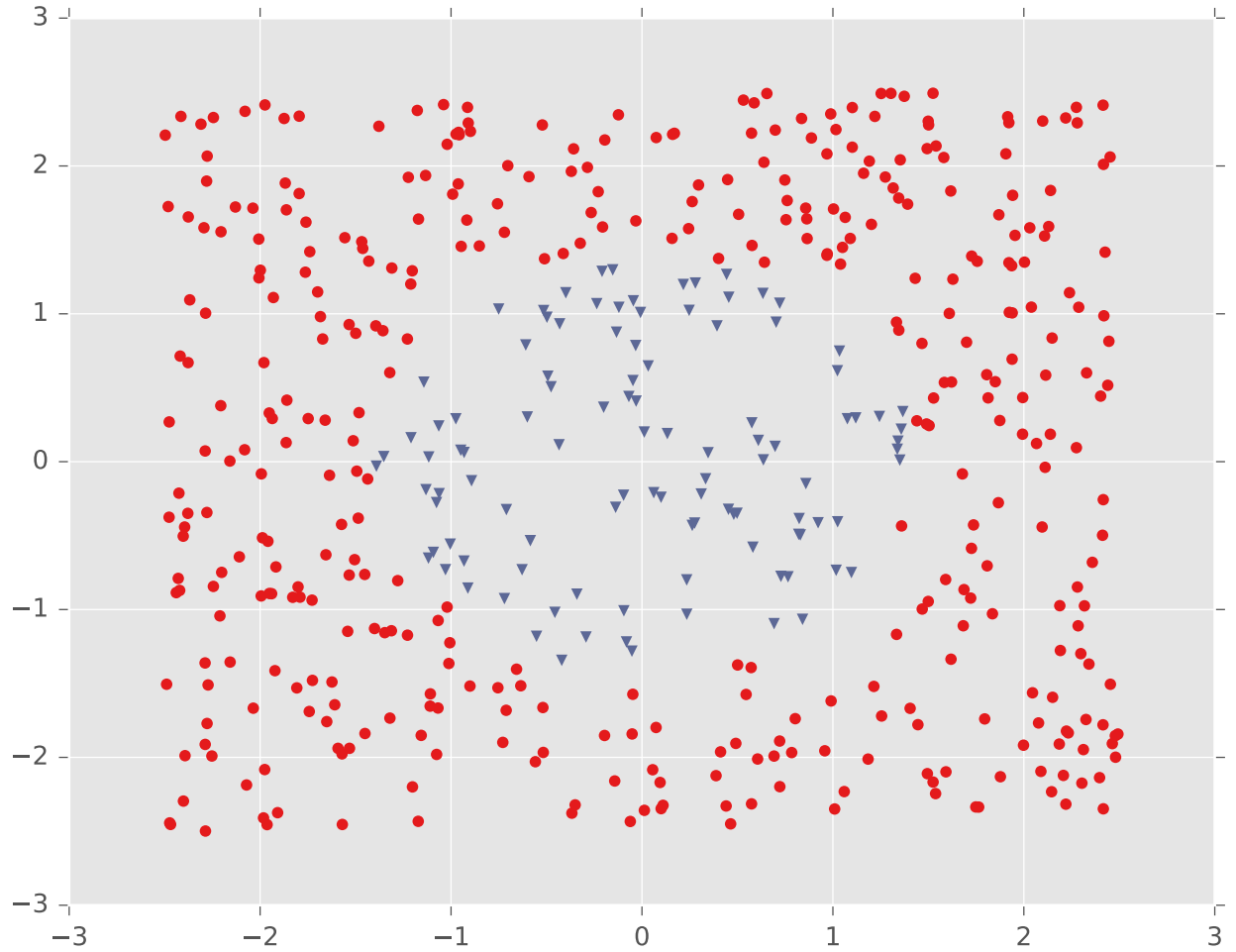
Visualizing
Gaussian
Naïve
Bayes
(3 classes,
equal
variances)



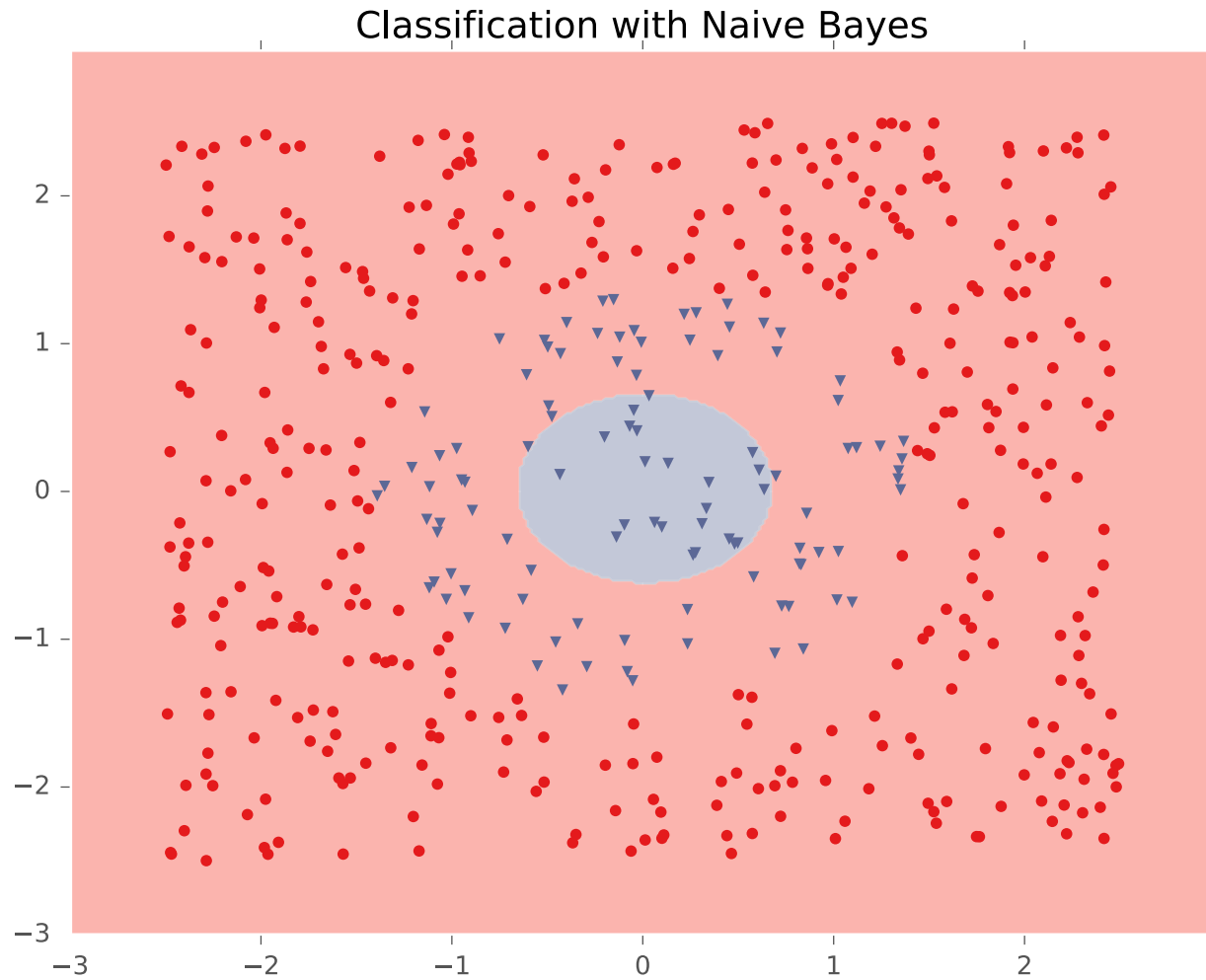
Visualizing Gaussian Naïve Bayes (3 classes, learned variances)



Visualizing Gaussian Naïve Bayes (2 classes, equal variances)



Visualizing
Gaussian
Naïve
Bayes
(2 classes,
~~equal
variances~~)



What if some
variables
never
appear in our
training data?
Predictions

- Given a test data point $\mathbf{x}' = [x'_1, \dots, x'_D]^T$

$$\hat{y} = \operatorname{argmax}_y \underbrace{P(Y=y | \mathbf{x}')}_{}$$

$$P(Y=1 | \mathbf{x}') \propto P(Y=1) P(\mathbf{x}' | Y=1)$$

$$\approx \hat{\pi} \left(\prod_{d=1}^D \hat{\theta}_{d,1}^{x'_d} (1 - \hat{\theta}_{d,1})^{1-x'_d} \right)$$

$$P(Y=0 | \mathbf{x}') \propto (1 - \hat{\pi}) \left(\prod_{d=1}^D \hat{\theta}_{d,0}^{x'_d} (1 - \hat{\theta}_{d,0})^{1-x'_d} \right)$$

$$\hat{y} = \begin{cases} 1 & \text{if } P(Y=1 | \mathbf{x}') \geq P(Y=0 | \mathbf{x}') \\ 0 & \text{otherwise} \end{cases}$$

What if some Word-Label pair never appears in our training data?

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

The Cat in the Hat gets a Dog (by ???)

- If some $\hat{\theta}_{d,y} = 0$ and that word appears in our test data \mathbf{x}' , then $P(Y = y|\mathbf{x}') = 0$ even if all the other features in \mathbf{x}' point to the label being y !

Setting the Parameters via MAP

- 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
 - $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$ and $\theta_{d,y} \sim \text{Beta}(\alpha, \beta)$
 - $\hat{\theta}_{d,y} = \frac{N_{Y=y, X_d=1} + (\alpha - 1)}{N_{Y=y} + (\alpha - 1) + (\beta - 1)}$
 - $N_{Y=y} = \#$ of data points with label y
 - $N_{Y=y, X_d=1} = \#$ of data points with label y and feature $X_d = 1$
- • Common choice: $\alpha = 2, \beta = 2$

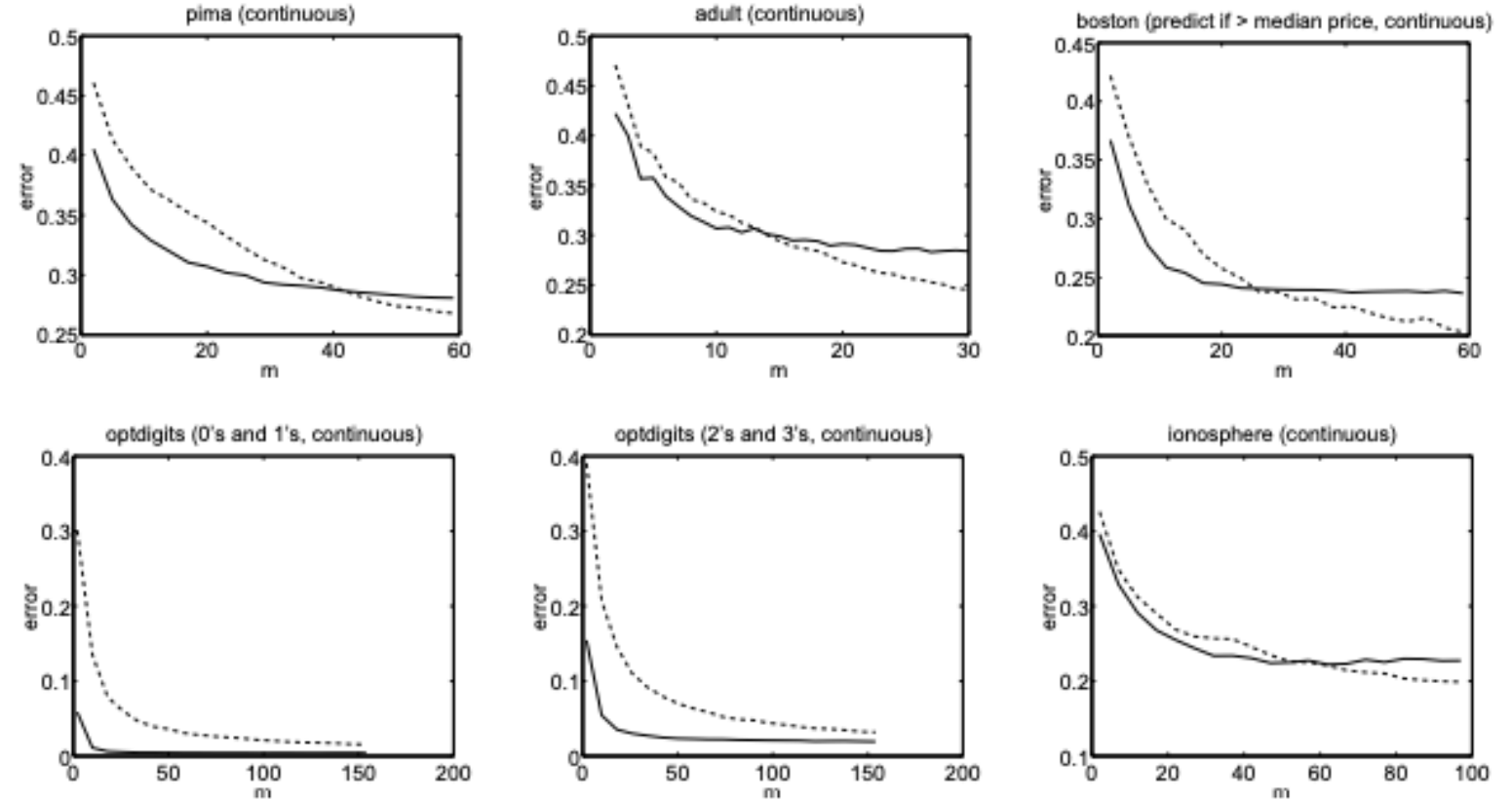
Logistic Regression vs. Naïve Bayes

- Naïve Bayes is a *generative* model
 - By modelling $P(X|Y)$ and $P(Y)$, we can *generate* new data points:
 1. Sample a label $y \sim P(Y)$
 2. Sample features $x_d \sim P(X_d|Y = y)$
- Logistic regression is a *discriminative* model
 - By modelling $P(Y|X)$, we can only *discriminate* (or distinguish) between classes.

Logistic Regression vs. Naïve Bayes (Ng and Jordan, 2001)

- Naïve Bayes and logistic regression form a *generative-discriminative* model pair
 - Recall that under certain conditions, the Gaussian Naïve Bayes (GNB) decision boundary is linear
 - If the Naïve Bayes assumption holds, then in the limit of infinite training data, GNB and logistic regression learn the same (linear) decision boundary!
 - In general, Naïve Bayes performs well when data is scarce but logistic regression has lower asymptotic error.

Logistic Regression vs. Naïve Bayes (Ng and Jordan, 2001)



- Dotted line: logistic regression
- Solid line: Naïve Bayes

Naïve Bayes Learning Objectives

You should be able to...

- Write the generative story for Naive Bayes
- Create a new naïve Bayes classifier using your favorite probability distribution as the event model
- Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of Bernoulli naïve Bayes
- Motivate the need for MAP estimation through the deficiencies of MLE
- Apply the principle of maximum a posteriori (MAP) estimation to learn the parameters of Bernoulli naïve Bayes
- Select a suitable prior for a model parameter
- Describe the tradeoffs of generative vs. discriminative models
- Implement Bernoulli naïve Bayes
- Describe how the variance affects whether a Gaussian naïve Bayes model will have a linear or nonlinear decision boundary