

ENGR-E 511; ENGR-E 399

Machine Learning for Signal Processing

Module 05: Bayesian Classification

Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

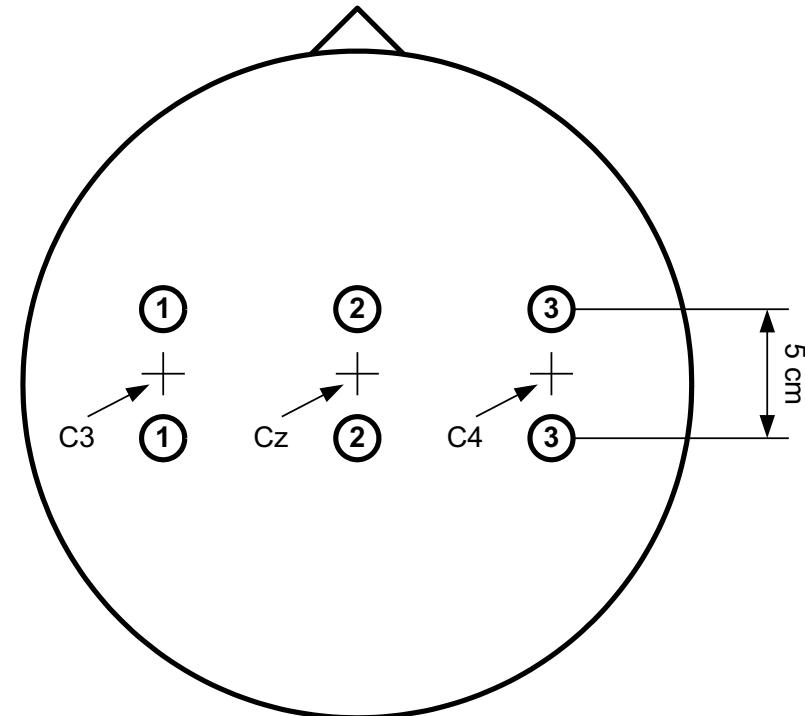
Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**



Brain Computer Interface

- <https://youtu.be/7t84IGE5TXA>



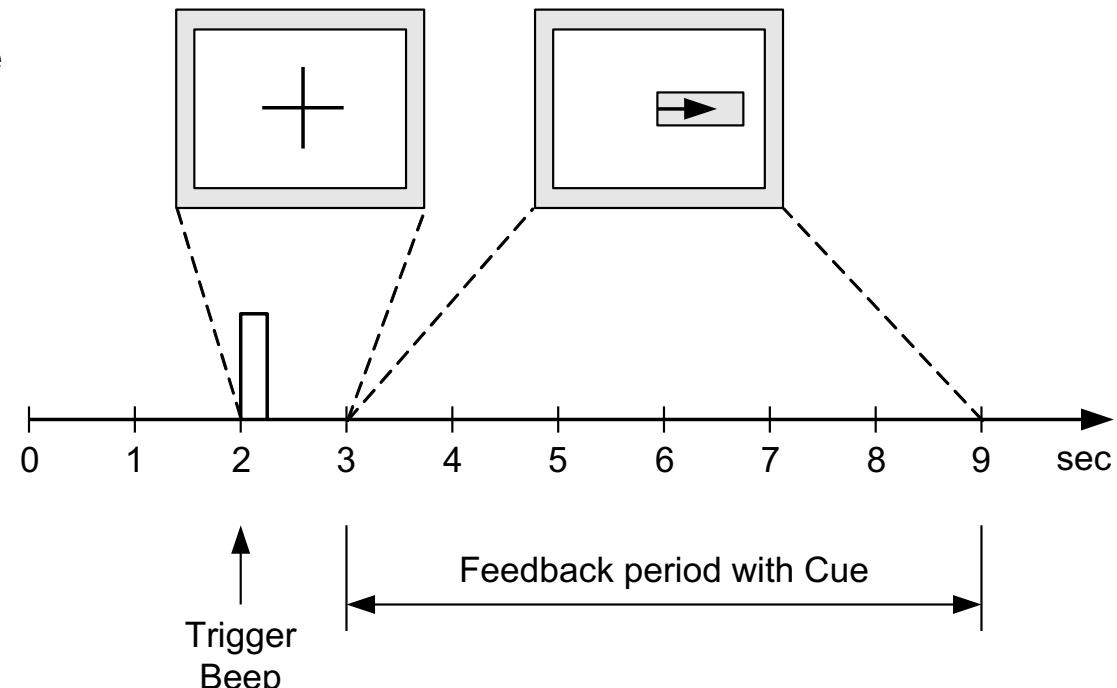
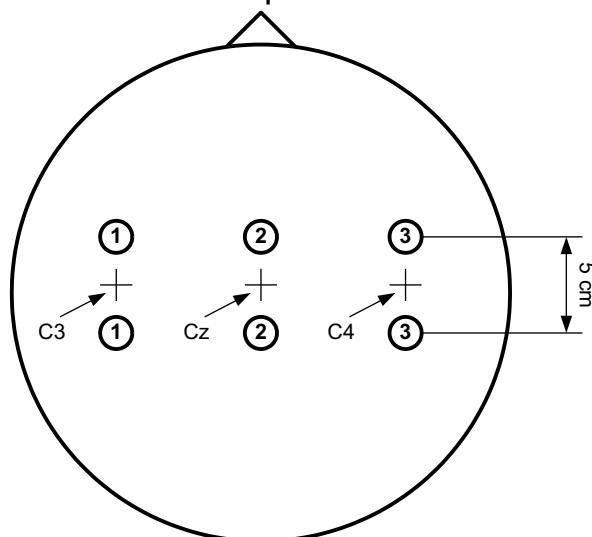
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

- Motor Imagery

- 112 recordings of 3-channel ElectroEncephaloGram (EEG)
- Each recording is 9-second long (128Hz)
 - After the triggers, from 3 to 9 sec the subject was asked to “imagine” the movement
 - Left or right
- If I keep only the last 6 sec after the trigger I have
 - 6 sec X 128 samples X 3 channels



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

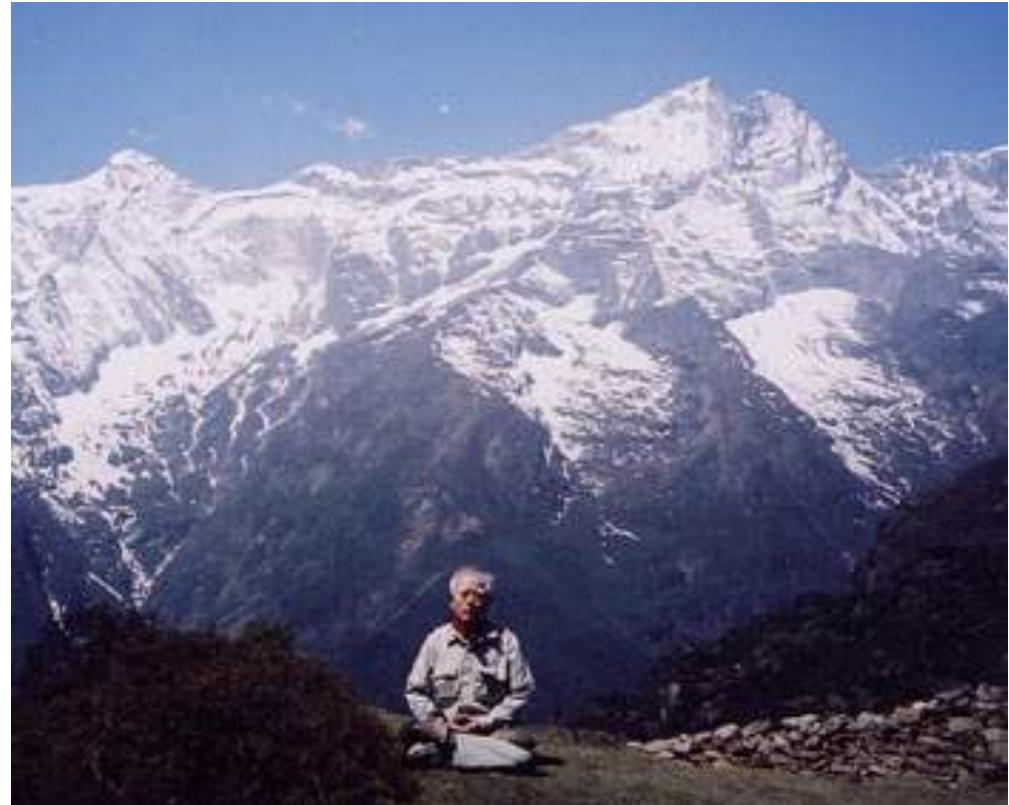
- Motor Imagery

- There are different brain waves with different frequencies

- Delta wave – (0.1 – 3 Hz)
- Theta wave – (4 – 7 Hz)
- Alpha wave – (8 – 13 Hz)
- Mu wave – (7.5 – 12.5 Hz)**
- SMR wave – (12.5 – 15.5 Hz)
- Beta wave – (16 – 31 Hz)
- Gamma wave – (32 – 100 Hz)



Meditation



- I care about Mu wave as it's known to be relevant to motor imagery
 - How to extract Mu wave?
 - STFT!



INDIANA UNIVERSITY

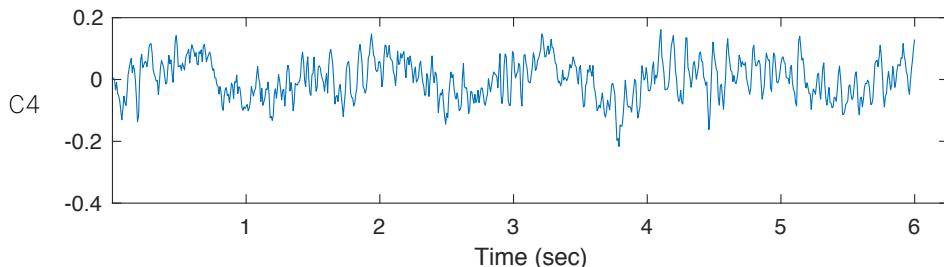
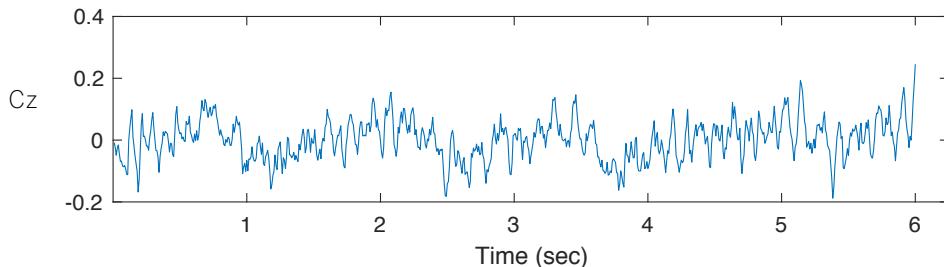
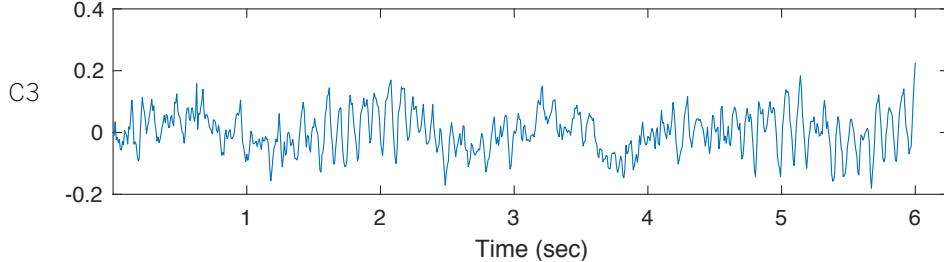
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

http://healthcare.joins.com/img/2003/0530/gra_0530_02.jpg

Brain Computer Interface

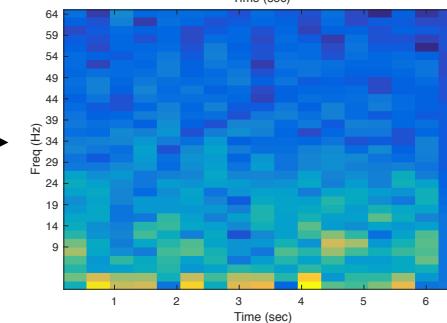
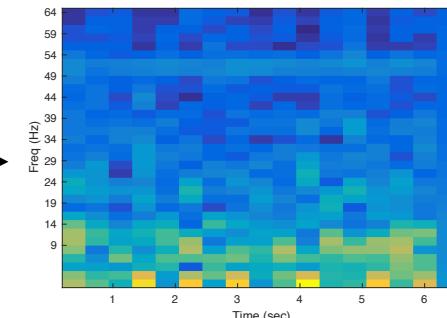
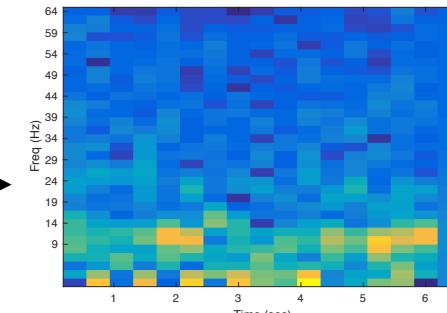
- Motor Imagery

○ STFT of the first trial



STFT
Frame size = 64
Hop size = 48
Window = "Blackman"

33 X 17



INDIANA UNIVERSITY

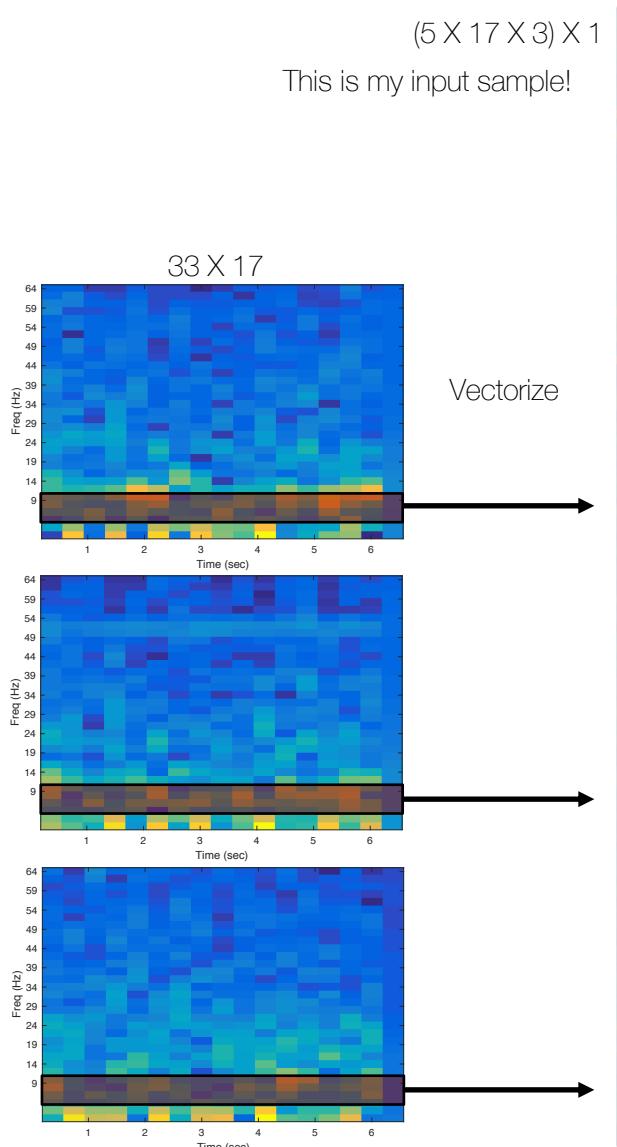
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

- Motor Imagery

- Because I care about the Mu wave – (7.5 – 12.5 Hz)
 - It's known to be relevant to motor imagery
- I only keep from 3rd to 7th rows
 - $64\text{Hz}/32 \times 3 = 6\text{Hz}$
 - $64\text{Hz}/32 \times 7 = 14\text{Hz}$
- Since I've got 112 trials like that
 - My data matrix is 255×112
- Now I'm going to use them to train a classifier
- What's wrong with my data set?
 - Too many dimensions!
 - Too small data (might be expensive)

(5 X 17 X 3) X 1
This is my input sample!



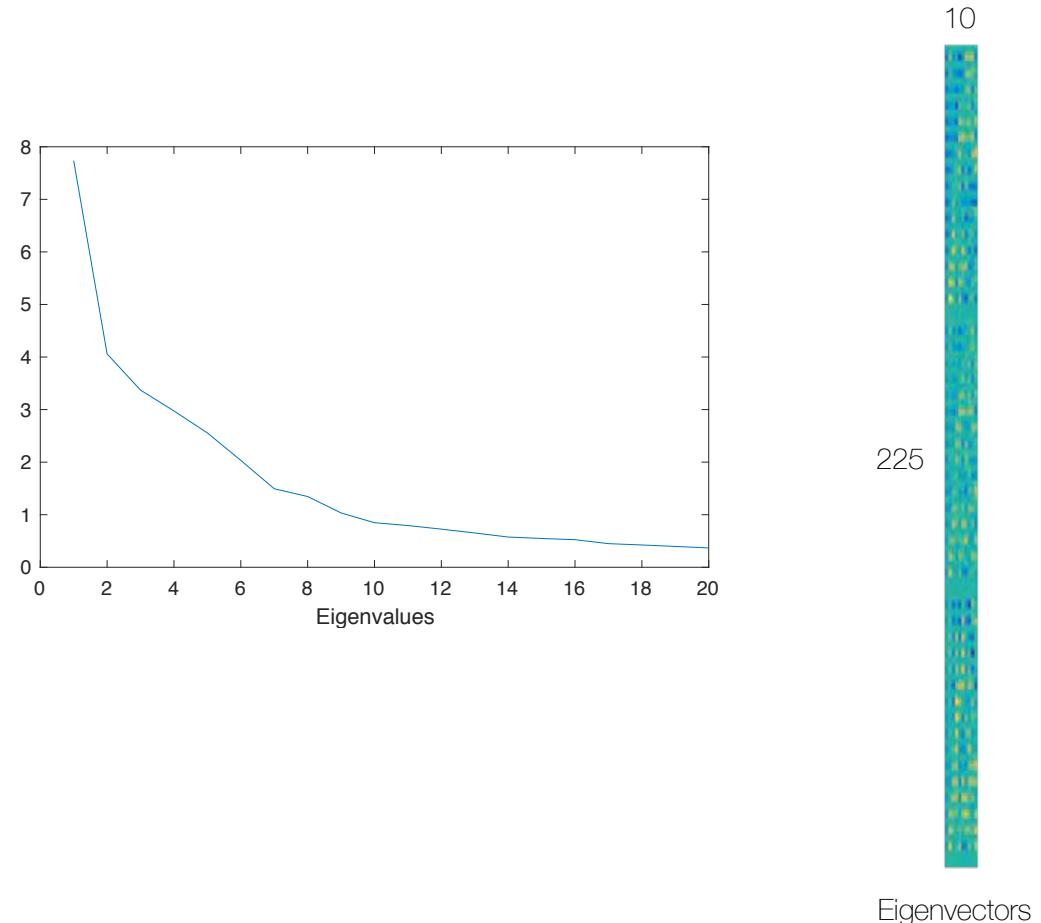
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

- Motor Imagery

- Now that we know how to reduce the dimension...
- I applied PCA on the data
- If you remember
 - Centering the data
 - Calculate the covariance matrix
 - Apply eigendecomposition
- I got a few eigenvectors and chose ten out of them
 - Why ten?
 - Based on the elbow of the eigenvalue graph



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

- Motor Imagery

- Now I can reduce the dimensionality from 225 down to 10

- How?

112

10

$Z_{(1:10,:)}$

255

$V^T_{(:,1:10)}$

112

255

X

- Let's see how they look like

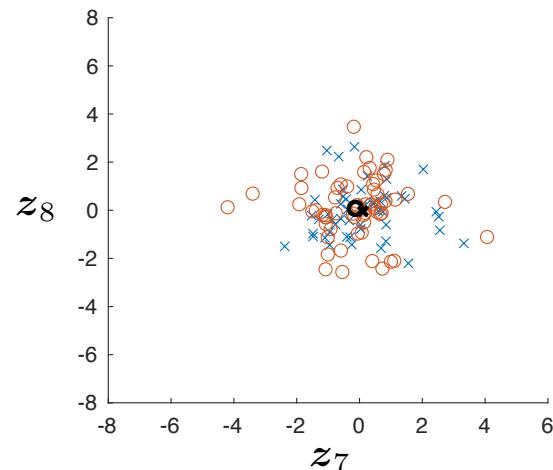
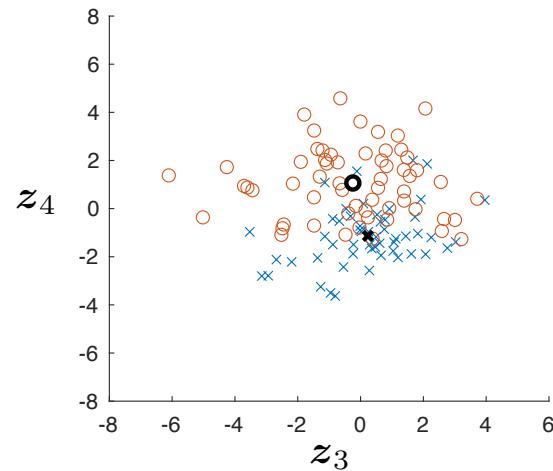
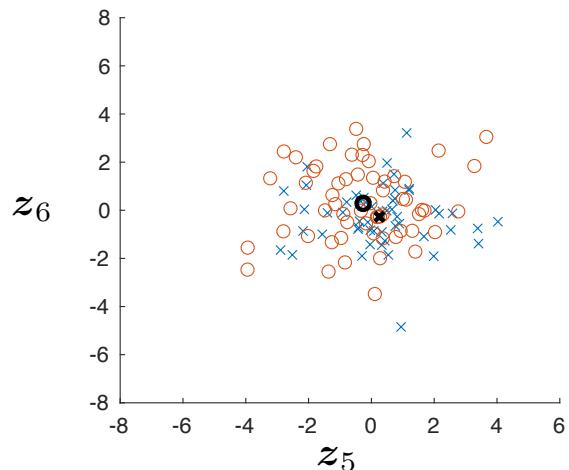
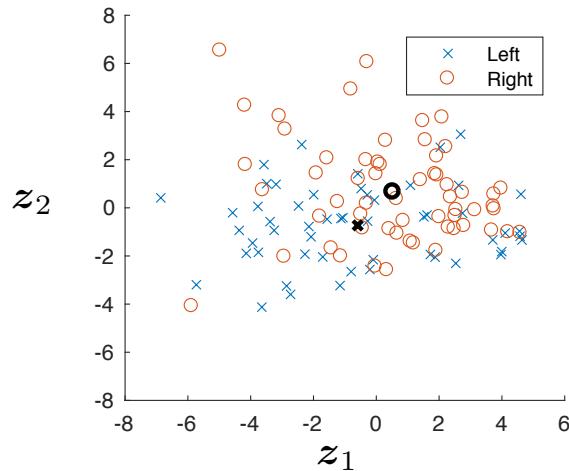


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Brain Computer Interface

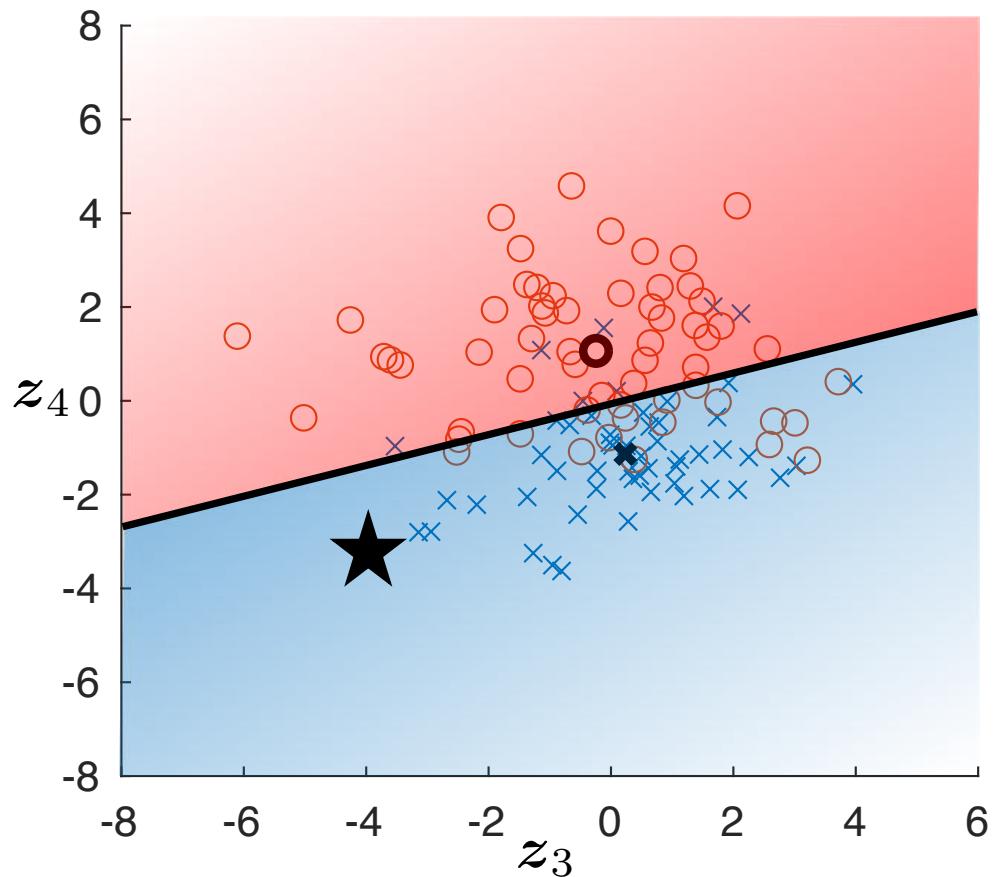
- Motor Imagery



Brain Computer Interface

- Motor Imagery

- For the convenience let's focus on z_3 and z_4
- I feel like drawing a line that can separate the two regions
 - If there's a test sample in the blue side, it might be from the "left" imagery
- How can I find the line?
 - There are many different ways!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

- I'll assume that each of the two classes follows a unique Gaussian dist.
- In other words, we need to *model* the training data
 - i.e. What kind of Gaussian should the underlying distribution be?
 - Any ideas?
- We do maximum likelihood estimation per class

$$\arg \max_{\boldsymbol{\mu}_L, \boldsymbol{\Sigma}_L} \sum_{i \in \mathcal{C}_L} \log \mathcal{N}(\mathbf{z}_{(:,i)}; \boldsymbol{\mu}_L, \boldsymbol{\Sigma}_L)$$

$$\arg \max_{\boldsymbol{\mu}_R, \boldsymbol{\Sigma}_R} \sum_{i \in \mathcal{C}_R} \log \mathcal{N}(\mathbf{z}_{(:,i)}; \boldsymbol{\mu}_R, \boldsymbol{\Sigma}_R)$$

- What would the estimation be?

$$\boldsymbol{\mu}_L^* = \frac{1}{|\mathcal{C}_L|} \sum_{i \in \mathcal{C}_L} \mathbf{z}_{(:,i)}$$

$$\boldsymbol{\mu}_R^* = \frac{1}{|\mathcal{C}_R|} \sum_{i \in \mathcal{C}_R} \mathbf{z}_{(:,i)}$$

$$\boldsymbol{\Sigma}_L^* = \frac{1}{|\mathcal{C}_L|} \sum_{i \in \mathcal{C}_L} (\mathbf{z}_{(:,i)} - \boldsymbol{\mu}_L^*)(\mathbf{z}_{(:,i)} - \boldsymbol{\mu}_L^*)^\top$$

$$\boldsymbol{\Sigma}_R^* = \frac{1}{|\mathcal{C}_R|} \sum_{i \in \mathcal{C}_R} (\mathbf{z}_{(:,i)} - \boldsymbol{\mu}_R^*)(\mathbf{z}_{(:,i)} - \boldsymbol{\mu}_R^*)^\top$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

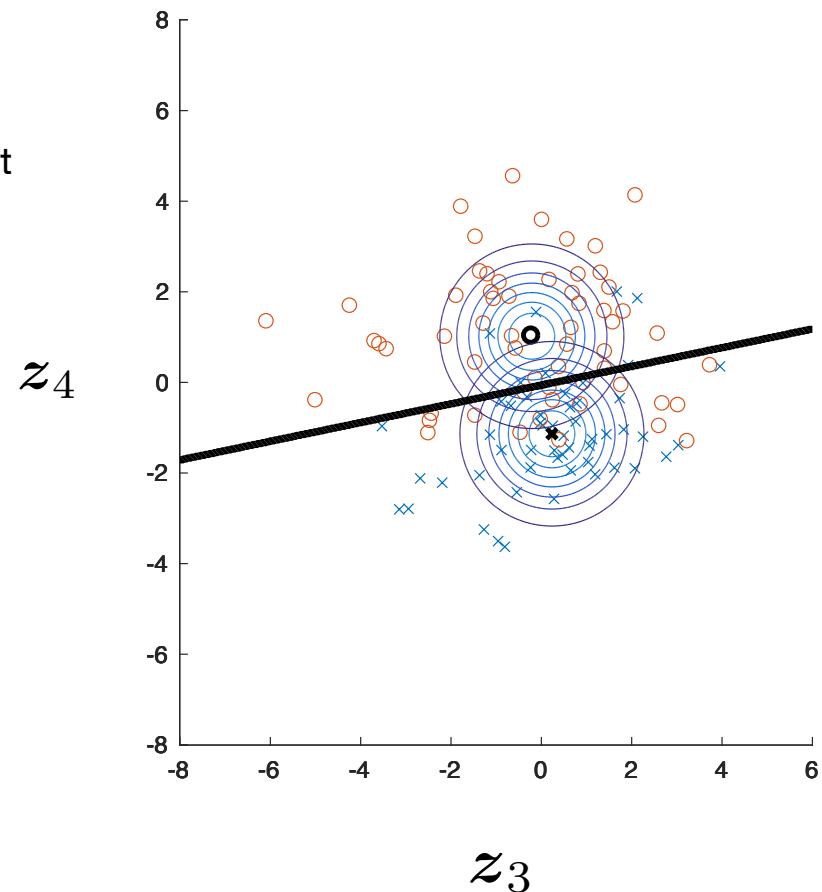
Naïve Bayes Classifier

- Decision Surface

- What would the hyperplane be?
- The Gaussians on the right are not optimal
 - Because I skipped the sample covariance matrix estimation part
 - Instead, I assumed that they are identity matrices

$$\Sigma_L^* = \Sigma_R^* = I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- Why would I do that?
 - Sometimes sample covariance matrix can be noisy
 - Sometimes? When?
 - If we don't have enough data samples
 - While their dimension is high



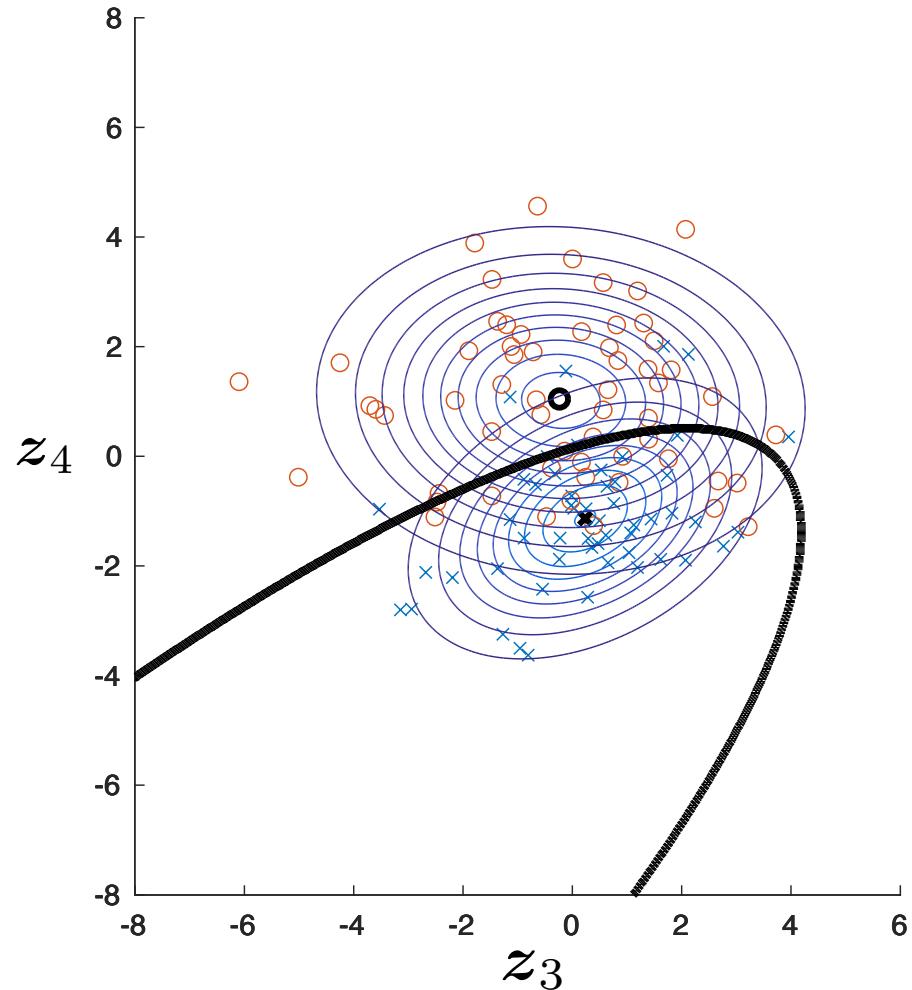
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

- Let's do it with the full covariance estimation
- Now we see that the decision surface is not linear
- Seems more flexible
 - More suitable for complicated classification problems
- Okay, modeling was easy,
but how do we find the hyperplane
 - And how do we make the decision
for the new unseen test samples?



INDIANA UNIVERSITY

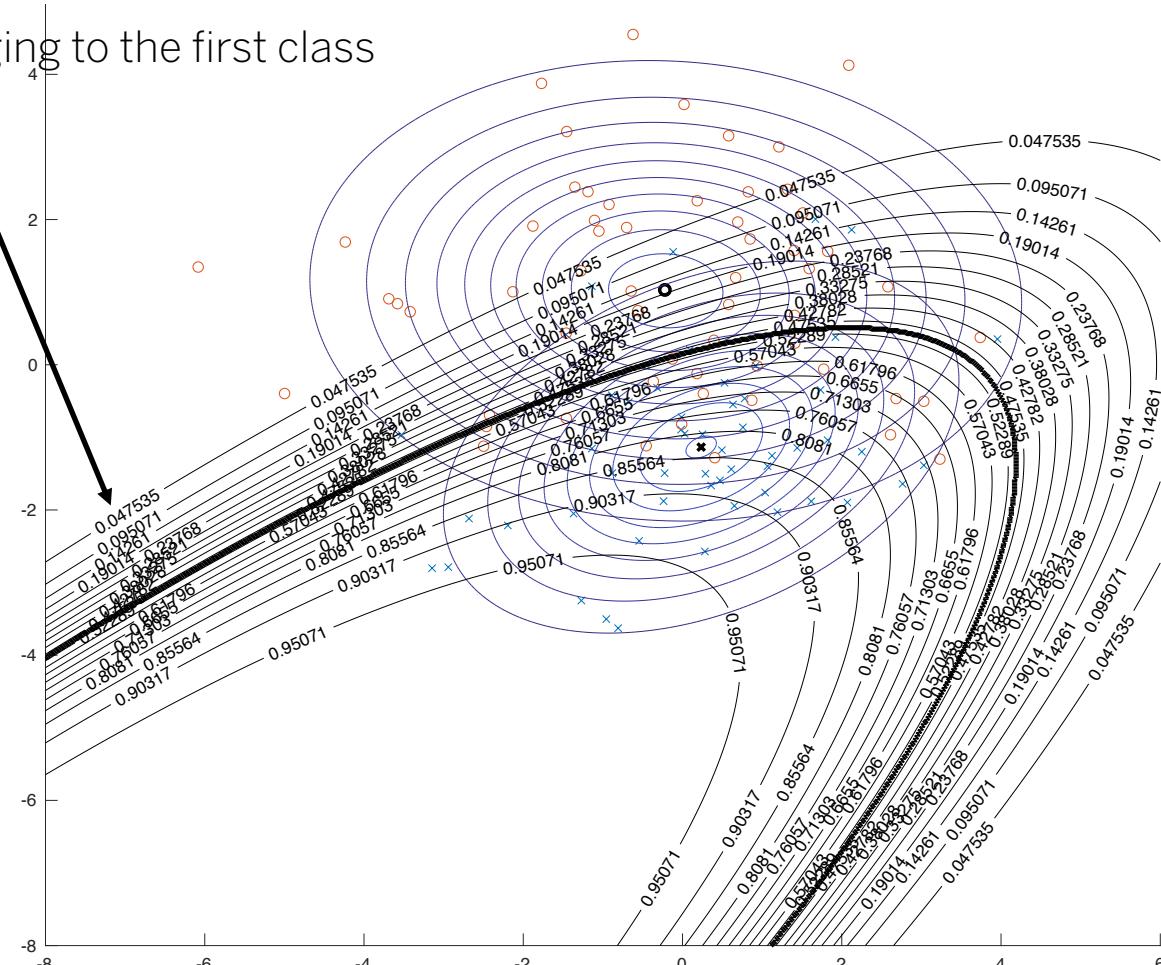
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

Probability of $z_{(:,i)}$ belonging to the first class

- Once again, how do we find these probabilities?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

- Now we need to use Bayes' rule
 - That's where the name is from
- We want to know the **posterior** probabilities of the classes given the data point

$$P(j|\mathbf{x}) = \frac{P(\mathbf{x}|j)P(j)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|j)P(j)}{\sum_j P(\mathbf{x}|j)P(j)}$$

- Or, for our EEG example

$$P(L|\mathbf{z}_{(:,i)}) = \frac{\mathcal{N}(\mathbf{z}_{(:,i)}; \boldsymbol{\mu}_L^*, \boldsymbol{\Sigma}_L^*)P(L)}{\mathcal{N}(\mathbf{z}_{(:,i)}; \boldsymbol{\mu}_L^*, \boldsymbol{\Sigma}_L^*)P(L) + \mathcal{N}(\mathbf{z}_{(:,i)}; \boldsymbol{\mu}_R^*, \boldsymbol{\Sigma}_R^*)P(R)}$$

↑ ↑
Prior probability of choosing "left" Prior probability of choosing "right"

$$P(j) = \frac{|\mathcal{C}_j|}{\sum_{j'=1}^J |\mathcal{C}_{j'}|}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

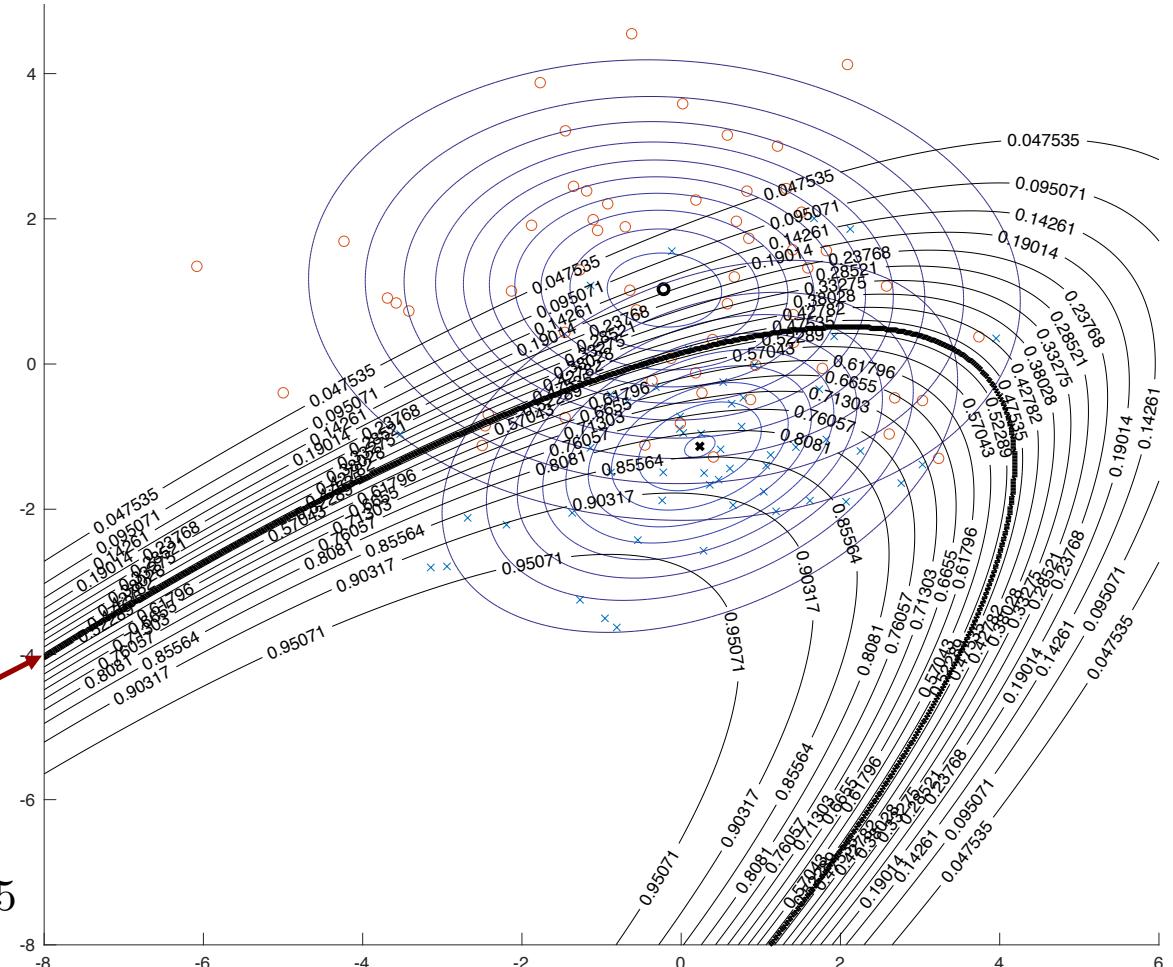
- Let's define the decision surface now

$$P(L|z_{(:,i)}) = P(R|z_{(:,i)})$$

- Or, if you like logarithm

$$\log P(L|z_{(:,i)}) = \log P(R|z_{(:,i)})$$

$$P(L|z_{(3,i)}, z_{(4,i)}) = P(R|z_{(3,i)}, z_{(4,i)}) = 0.5$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface

- In general for the naïve Bayes, the decision boundary between the class j and k meets this condition:

$$\log P(j|x) = \log P(k|x)$$

$$\log P(j)\mathcal{N}(x; \mu_j, \Sigma_j) = \log P(k)\mathcal{N}(x; \mu_k, \Sigma_k)$$

$$\begin{aligned} & \log P(j) + (\mathbf{x} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) + const. \\ &= \log P(k) + (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + const. \end{aligned}$$

$$\begin{aligned} & \log P(j) + \cancel{\mathbf{x}^\top \boldsymbol{\Sigma}_j^{-1} \mathbf{x}} + 2\mathbf{x}^\top \boldsymbol{\Sigma}_j^{-1} \boldsymbol{\mu}_j + \boldsymbol{\mu}_j^\top \boldsymbol{\Sigma}_j^{-1} \boldsymbol{\mu}_j \\ &= \log P(k) + \cancel{\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{x}} + 2\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \end{aligned}$$

- When is the decision boundary linear?
 - When there's no quadratic terms in the equation
 - i.e. when the two covariance matrices are same

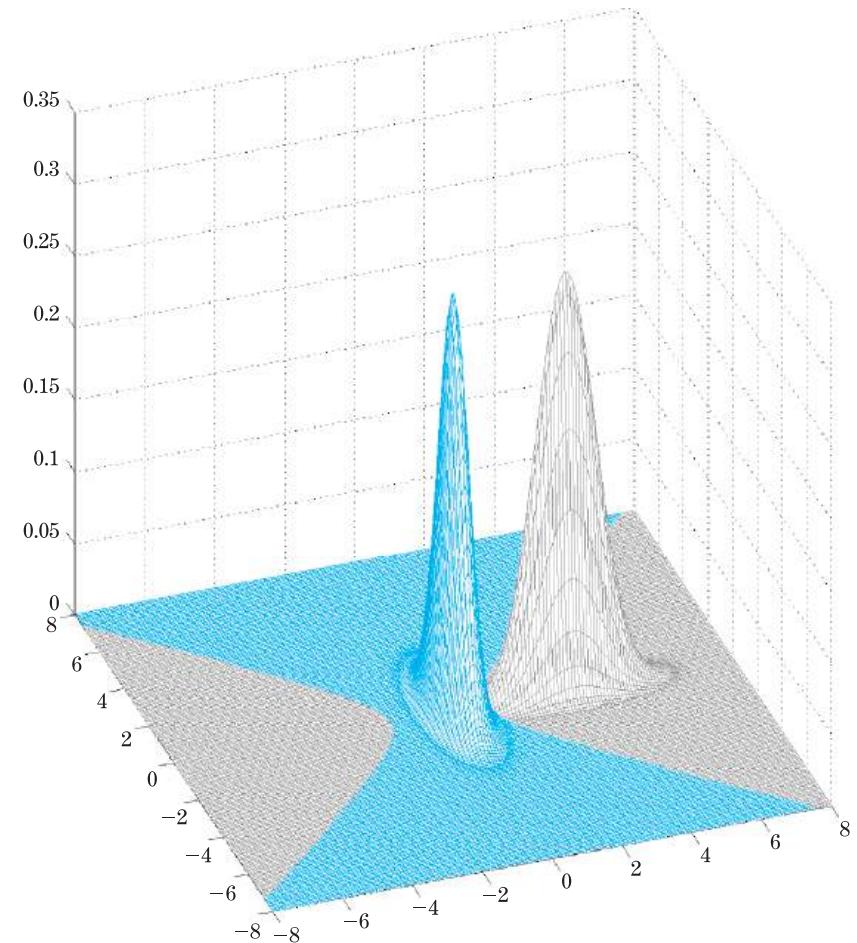
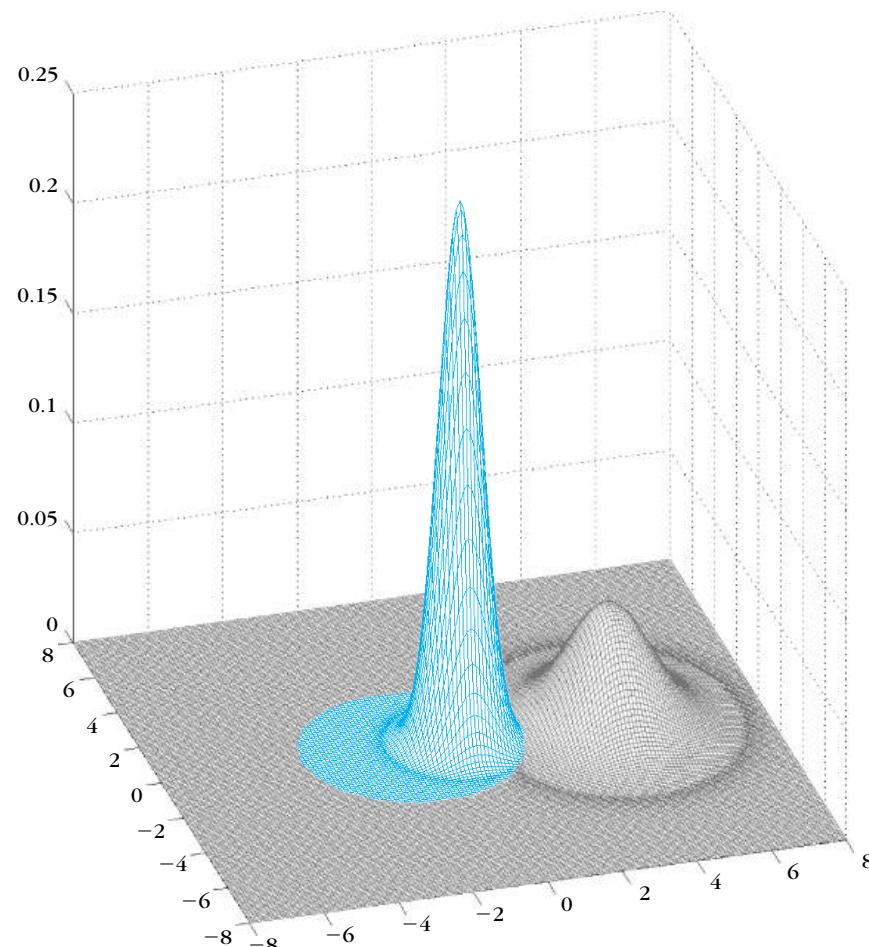


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Naïve Bayes Classifier

- Decision Surface



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

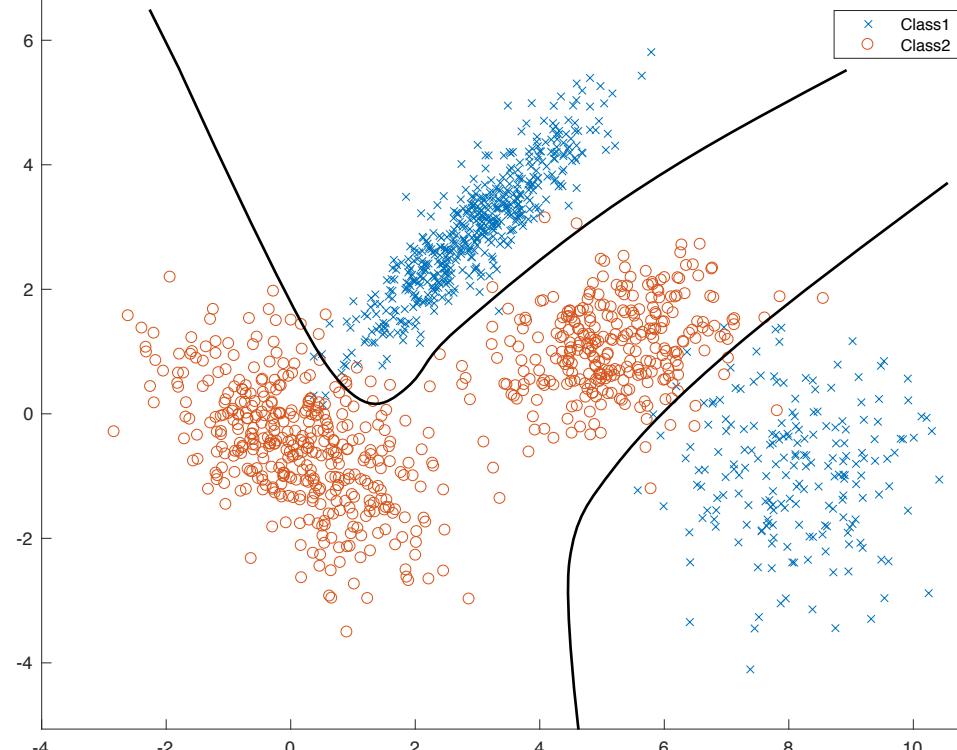
Figure 2.8 and 2.9

GMM Classifier

- Not for clustering, but for classification!

- What are we going to do with this?

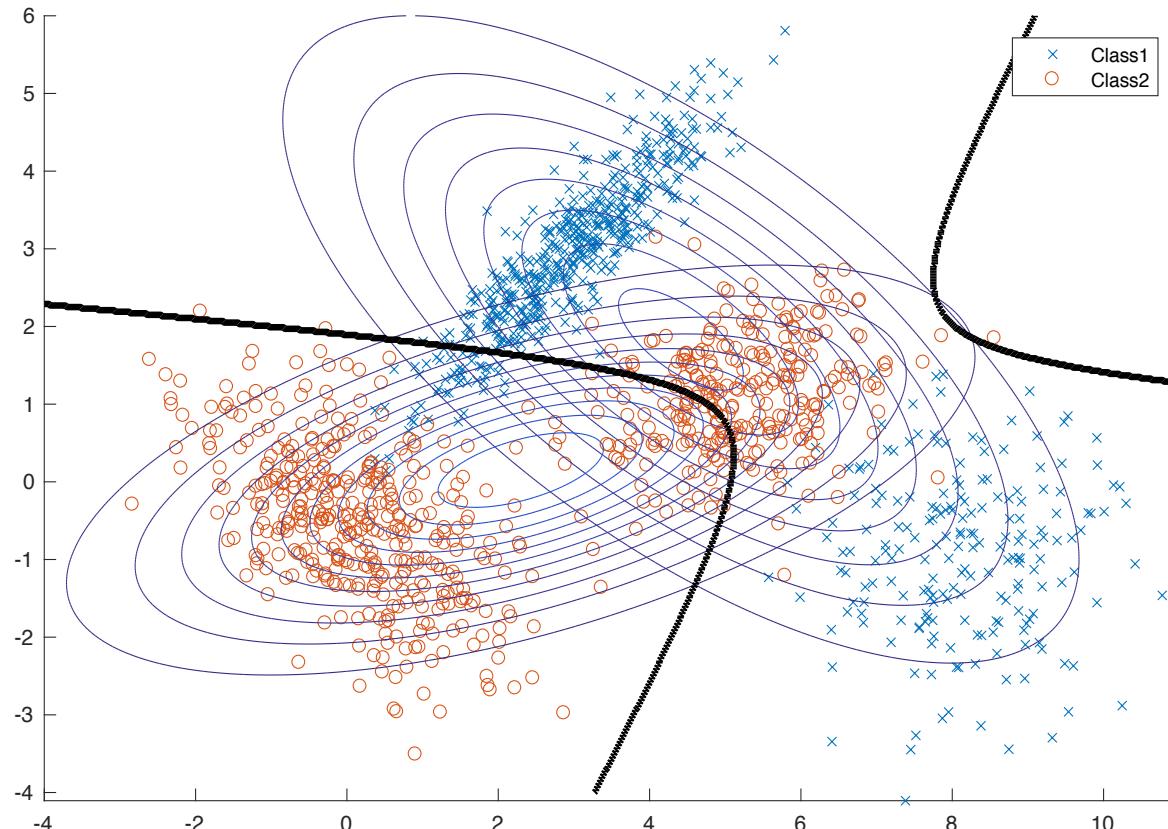
- I feel like drawing a nonlinear decision boundary
 - Can we do this with naïve Bayes?



GMM Classifier

- Not for clustering, but for classification!

- Naïve Bayes results are not what we want

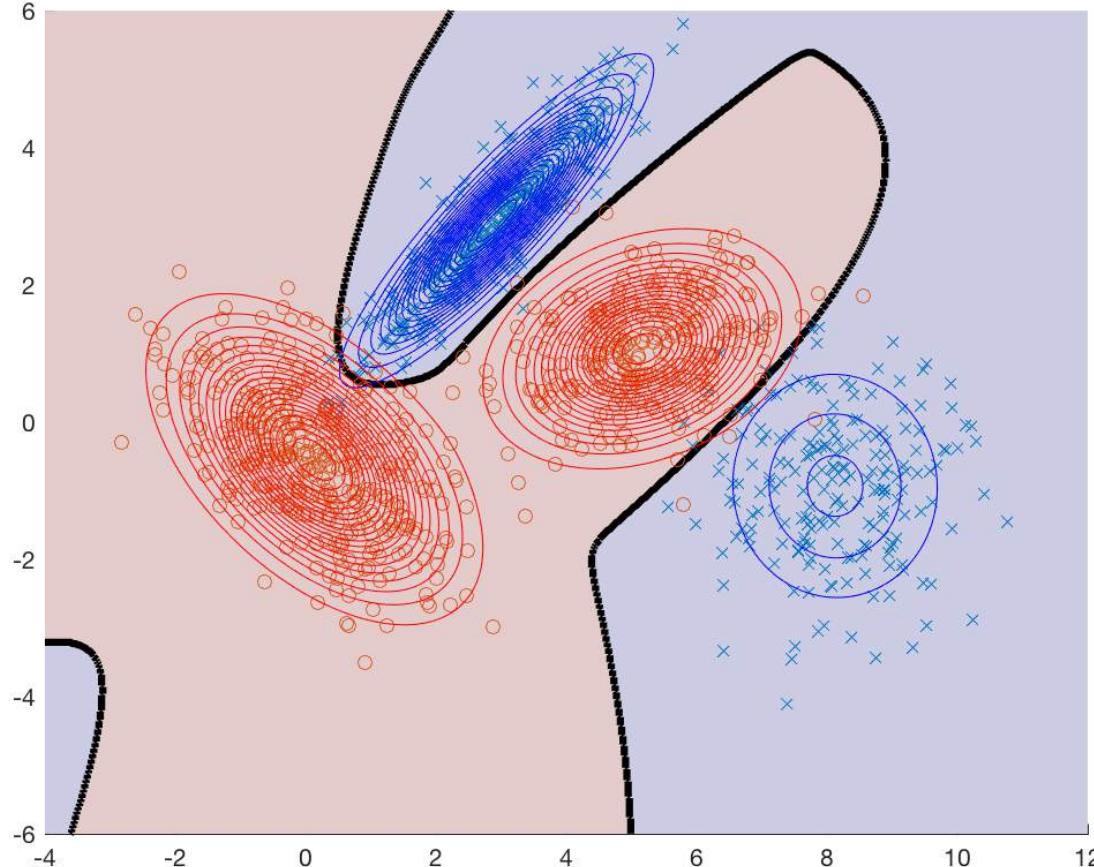


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

GMM Classifier

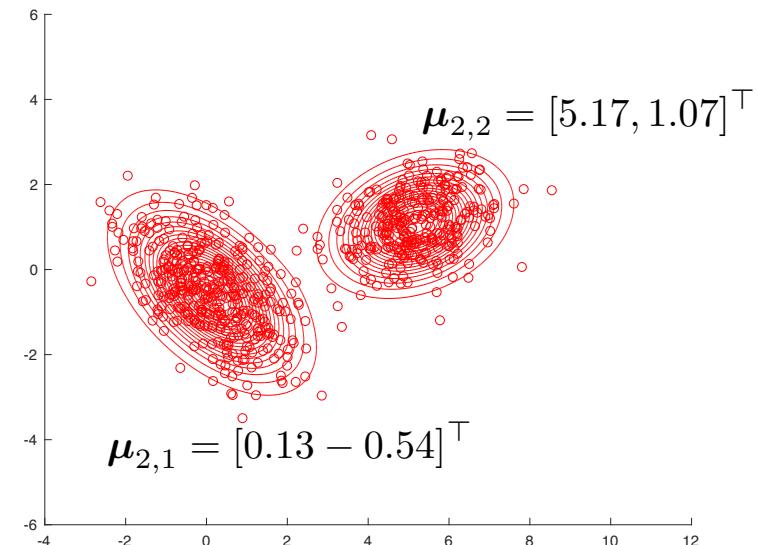
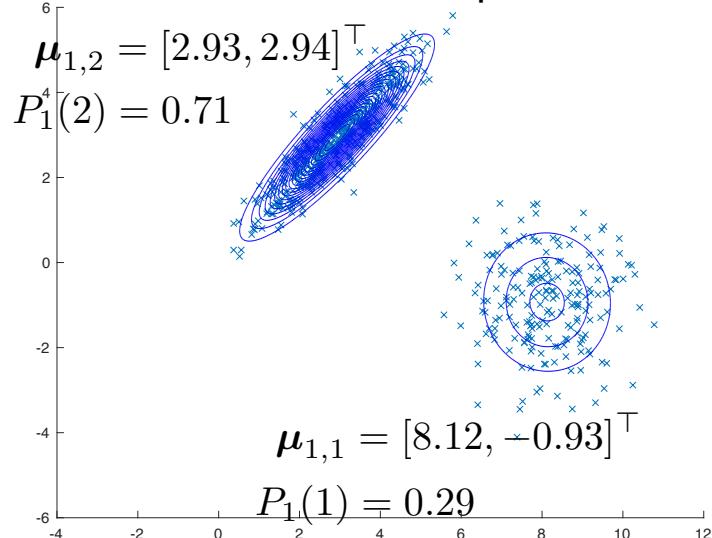
- Not for clustering, but for classification!
- GMM can define highly non-linear decision boundaries



GMM Classifier

- Not for clustering, but for classification!

- Math?
- First we need to learn a GMM per class



- This will give me J sets of parameters
 - J is the number of classes; L_j is the number of Gaussians in the j -th class
- $$j = 1 : \{(\mu_{1,1}, \Sigma_{1,1}), P_1(1)\}, \{(\mu_{1,2}, \Sigma_{1,2}), P_1(2)\}, \dots, \{(\mu_{1,L_1}, \Sigma_{1,L_1}), P_1(L_1)\}$$
- $$j = 2 : \{(\mu_{2,1}, \Sigma_{2,1}), P_2(1)\}, \{(\mu_{2,2}, \Sigma_{2,2}), P_2(2)\}, \dots, \{(\mu_{2,L_2}, \Sigma_{2,L_2}), P_2(L_2)\}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

GMM Classifier

- Not for clustering, but for classification!

- During the test time, since we know all the parameters we need
- We're going to apply the same Bayes rule

$$P(j|\mathbf{x}) = \frac{P(\mathbf{x}|j)P(j)}{\sum_j P(\mathbf{x}|j)P(j)}$$

- Some additional summations to take care of GMM inside the classifier

$$\text{i.e. } P(j) = \frac{|\mathcal{C}_j|}{\sum_{j'=1}^J |\mathcal{C}_{j'}|}$$

Prior for the j -th class

The summation over l_j is for this j -th class only

GMM for the j -th class

$$P(j|\mathbf{x}) = \frac{P(j) \sum_{l_j=1}^{L_j} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{j,l_j}, \boldsymbol{\Sigma}_{j,l_j}) P_j(l_j)}{\sum_j P(j) \sum_{l_j=1}^{L_j} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{j,l_j}, \boldsymbol{\Sigma}_{j,l_j}) P_j(l_j)}$$



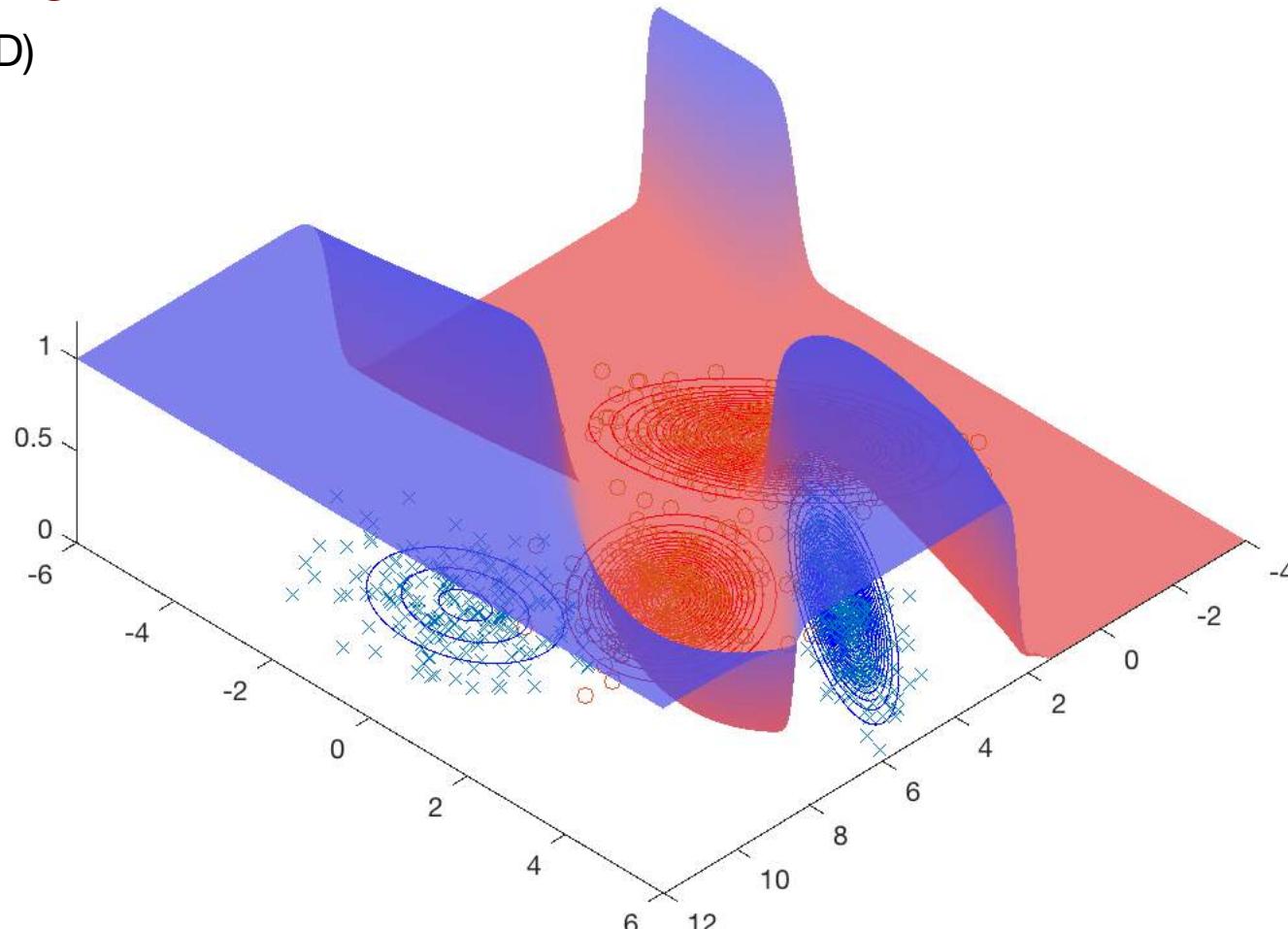
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

GMM Classifier

- Not for clustering, but for classification!

- GMM classifier (3D)



Non-parametric Classifier

- Parzen-Rosenblatt windows for kernel density estimation: let's go finer
 - What if we assume much more Gaussians per class?
 - Like, $L_j = |\mathcal{C}_j|$ Gaussians?

- Too many parameters to estimate during training

- No way to calculate the covariance matrix

- Let's forget about the covariance matrices

$$\Sigma_{j,l} = \sigma^2 \mathbf{I}, \quad \forall j, l$$

- Let's forget about the means: assume mean is the training samples

$$\mu_{j,l} = \mathbf{x}_i, \quad \text{where } i \in \{1, 2, \dots, L_j\}$$



INDIANA UNIVERSITY

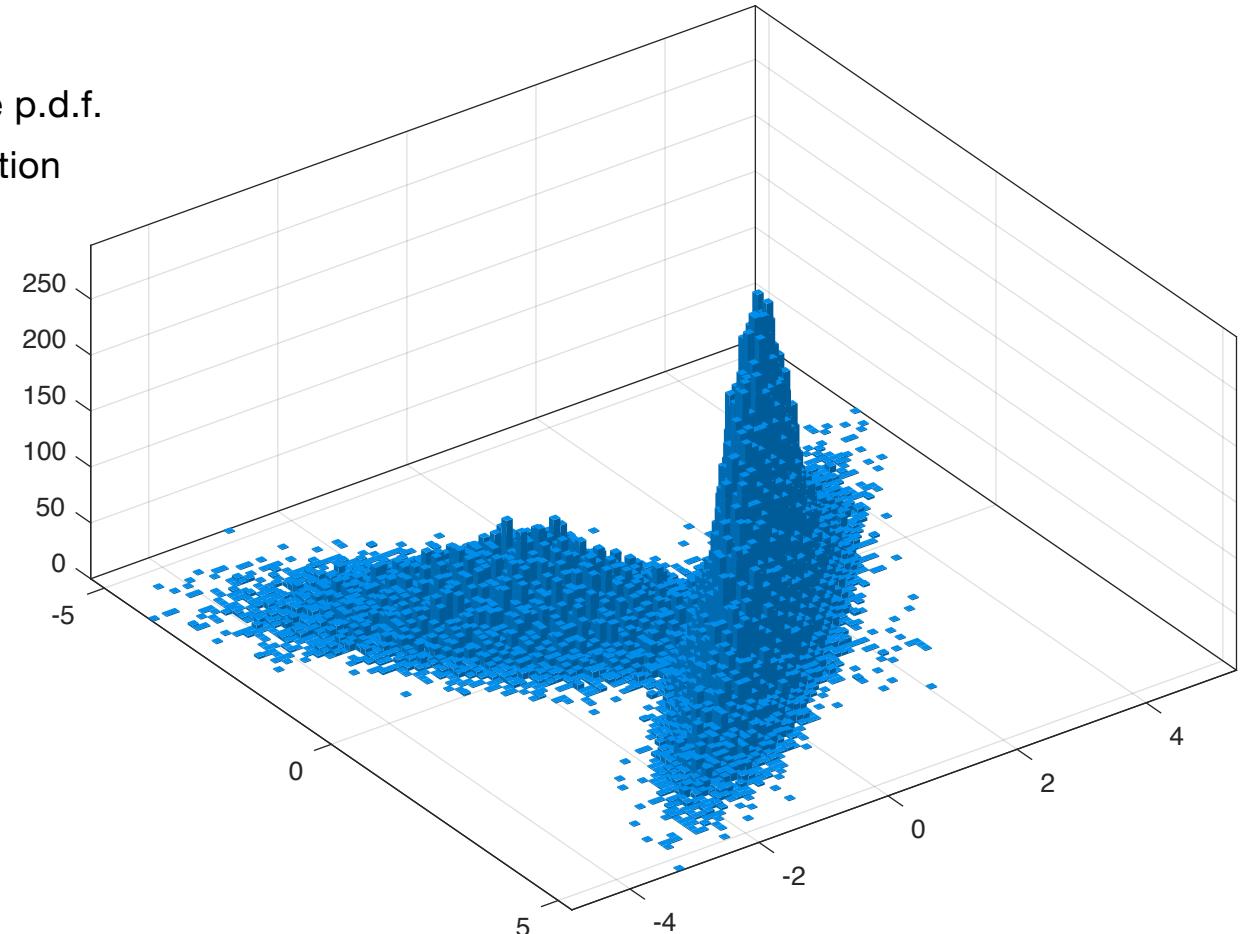
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- Parzen-Rosenblatt windows for kernel density estimation: let's go finer

- Non-parametric density estimation

- Doesn't assume any parametric function for the p.d.f.
 - Data says something about the sample distribution
 - e.g. Histograms



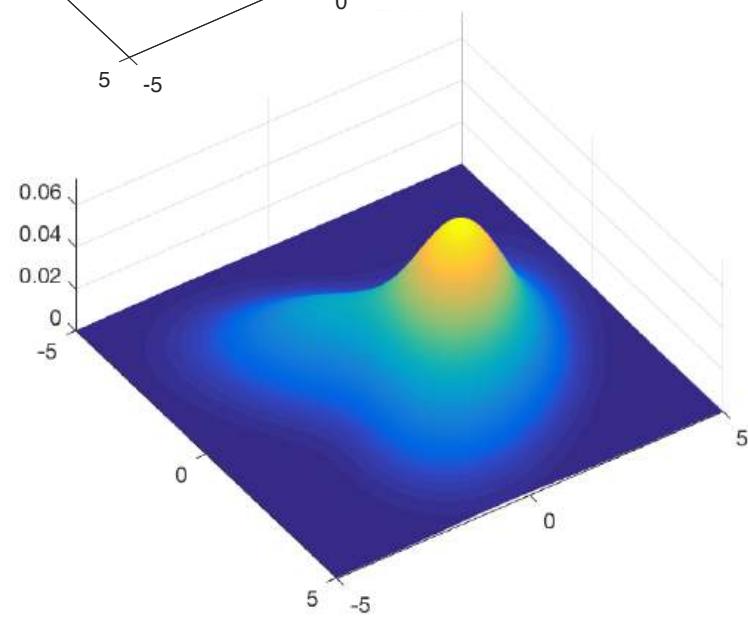
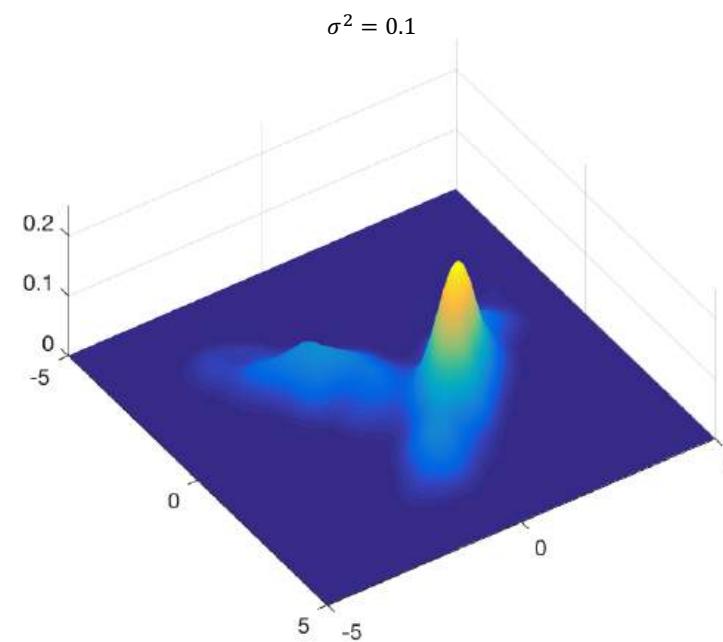
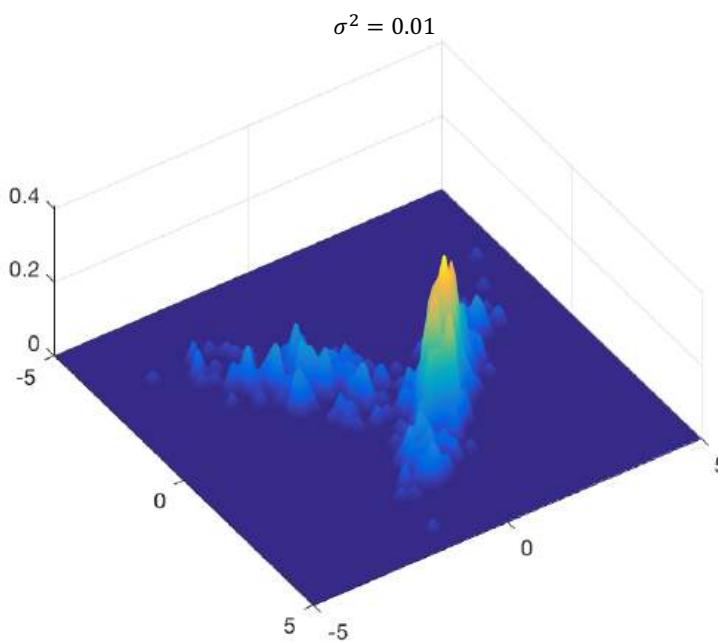
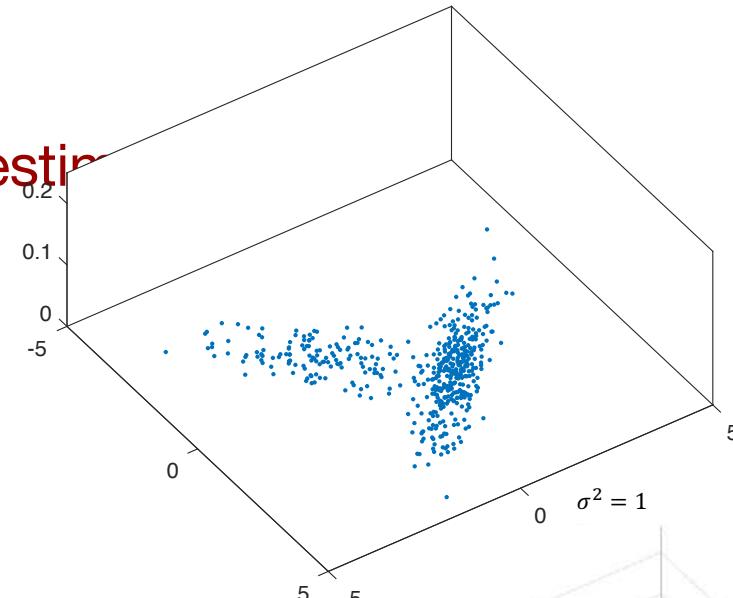
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- Parzen-Rosenblatt windows for kernel density estimation

- In kernel density estimation
 - Each sample forms a small Gaussian with an isotropic covariance



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- Parzen-Rosenblatt windows for kernel density estimation: let's go finer

○ In k-means clustering what we assume about the covariance is

- Isotropic (spherical) $\sigma_j \mathbf{I}$
- $\sigma_1 = \sigma_2 = \dots = \sigma_J$

○ On top of clustering, another purpose of GMM was

- Density estimation

○ KDE using Gaussian kernels

$$p(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(\mathbf{x}; \mathbf{x}_i, \sigma^2 \mathbf{I})$$

○ In general

$$p(\mathbf{x}) = \frac{1}{N\sigma} \sum_{i=1}^N \mathcal{K}\left(\frac{\mathbf{x} - \mathbf{x}_i}{\sigma}\right)$$

- Integral of the kernel function \mathcal{K} should be 1

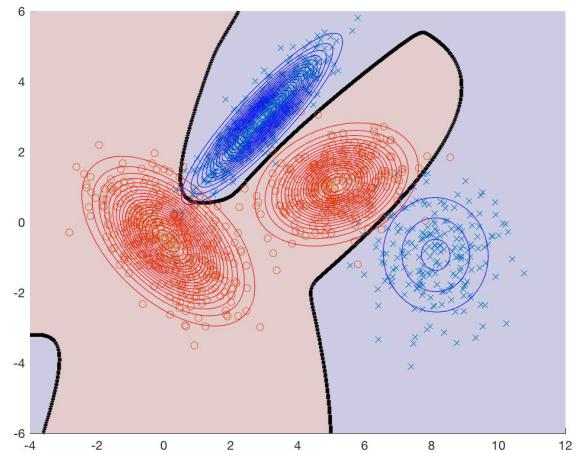


INDIANA UNIVERSITY

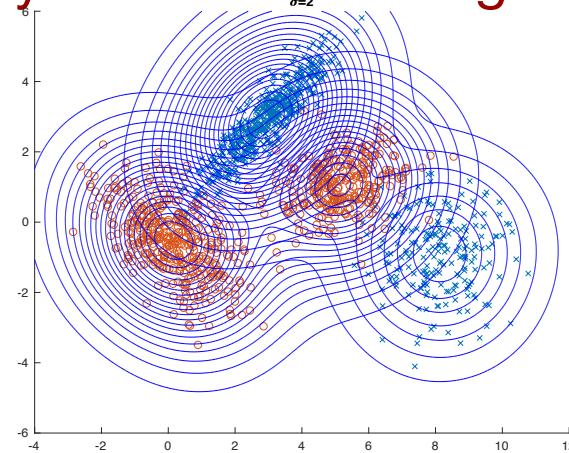
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

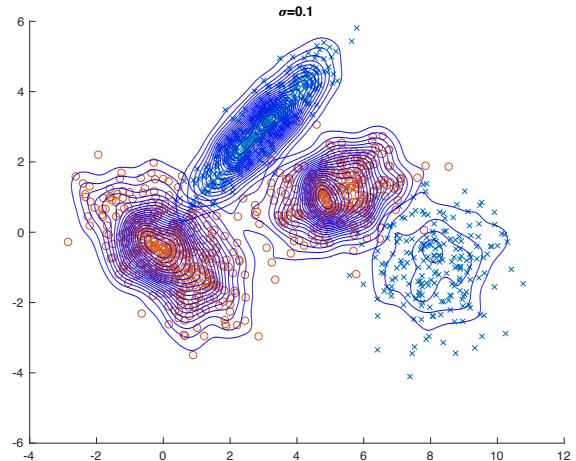
- Parzen-Rosenblatt windows for kernel density estimation: let's go finer



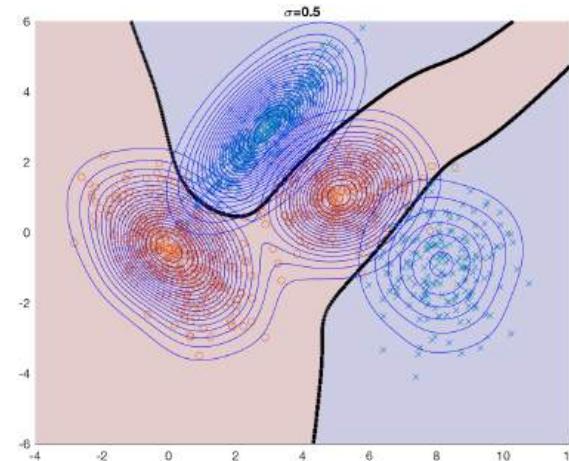
GMM Classifier



KDE Classifier
 $\sigma^2 = 2$



KDE Classifier
 $\sigma^2 = 0.1$



KDE Classifier
 $\sigma^2 = 0.5$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- Parzen-Rosenblatt windows for kernel density estimation: let's go finer
 - I don't like it, because
 - Too much fuss
 - For a test sample, its posterior probability involves all training samples

$$\begin{aligned} P(j|x) &= \frac{\frac{|\mathcal{C}_j|}{N} \sum_{i \in \mathcal{C}_j} \mathcal{N}(x; \mathbf{x}_i, \sigma^2 \mathbf{I}) / |\mathcal{C}_j|}{\sum_j \frac{|\mathcal{C}_j|}{N} \sum_{i \in \mathcal{C}_j} \mathcal{N}(x; \mathbf{x}_i, \sigma^2 \mathbf{I}) / |\mathcal{C}_j|} \\ &= \frac{\sum_{i \in \mathcal{C}_j} \mathcal{N}(x; \mathbf{x}_i, \sigma^2 \mathbf{I})}{\sum_j \sum_{i \in \mathcal{C}_j} \mathcal{N}(x; \mathbf{x}_i, \sigma^2 \mathbf{I})} \end{aligned}$$

- [Sum of all the pdf defined by all the training samples of that class] / [Sum of all the pdf defined by all the training samples]
 - Not a very attractive feature as a classifier (e.g. for a real-time application)



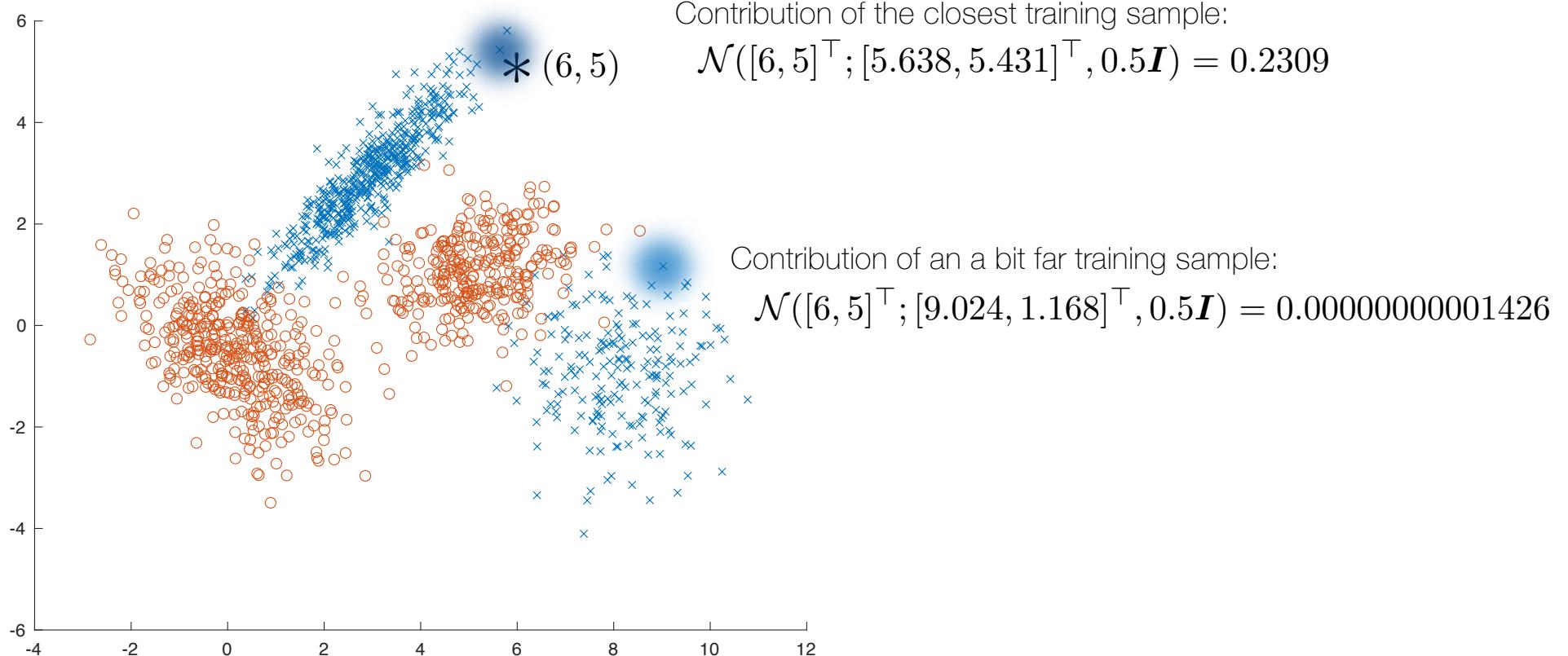
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- k-Nearest Neighbors

- Why do we care about the small Gaussian distribution that's far away from my current test sample?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- k-Nearest Neighbors

- kNN density estimation (for j -th class)
 - We assume a region \mathcal{R} tightly defined with the k neighbors
 - Its volume V depends on the density of the region
 - This \mathcal{R} is our kernel, whose pdf is uniform (not Gaussian)

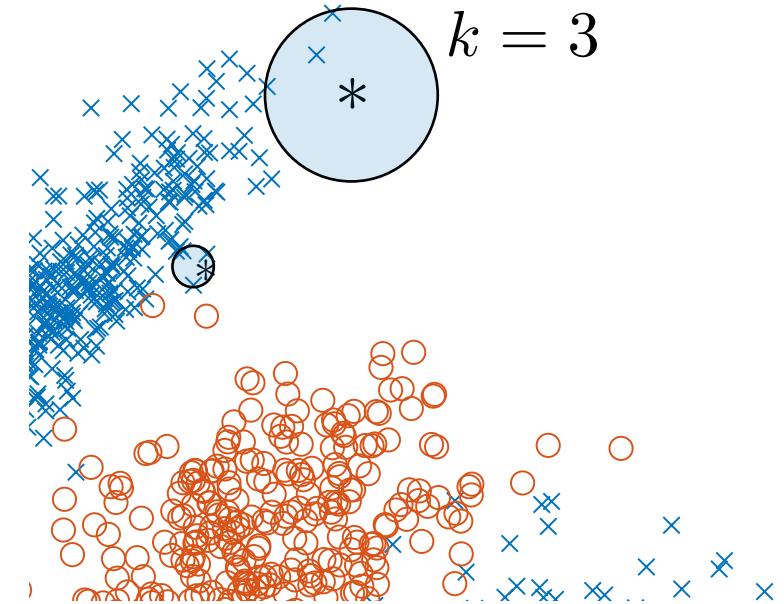
$$P(\mathbf{x}|l) = \frac{1}{V}$$

- What's the prior probability?
 - Probability of choosing this particular kernel?

$$P(l) = \frac{k}{N}$$

- Therefore, the density of j -th class

$$\begin{aligned} P(\mathbf{x}|j) &= \sum_{l=1}^{L_j} P(\mathbf{x}|l)P(l) \\ &= P(\mathbf{x}|l)P(l) \leftarrow \end{aligned}$$



because l is the only kernel \mathbf{x} belongs to according to the kNN definition



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- k-Nearest Neighbors

- k NN classifier

- Once again, the same Bayes rule

$$P(j|\mathbf{x}) = \frac{P(j) \sum_{l_j=1}^{L_j} P(\mathbf{x}|l_j)P(l_j)}{\sum_{j=1}^J P(j) \sum_{l_j=1}^{L_j} P(\mathbf{x}|l_j)P(l_j)} = \frac{P(j)P(\mathbf{x}|l_j)P(l_j)}{\sum_{j=1}^J P(j)P(\mathbf{x}|l_j)P(l_j)}$$

- But this time we see only one kernel, formed with the k NN
 - Plug in the k NN probabilities
- $$P(\mathbf{x}|l_j) = \frac{1}{V_{l,j}} \quad P(l_j) = \frac{k}{|\mathcal{C}_j|} \quad P(j) = \frac{|\mathcal{C}_j|}{N}$$
- Then, $P(j|\mathbf{x}) = \frac{\frac{|\mathcal{C}_j|}{N} \frac{k}{V_{l,j}|\mathcal{C}_j|}}{P(\mathbf{x})} = \frac{k}{V_{l,j}NP(\mathbf{x})} = \frac{Const}{V_{l,j}}$
 - **If the volume of kernel is inverse-proportional to the posterior probability**
- $$V_{l,j} > V_{l,j'} \rightarrow P(j'|\mathbf{x}) > P(j|\mathbf{x})$$



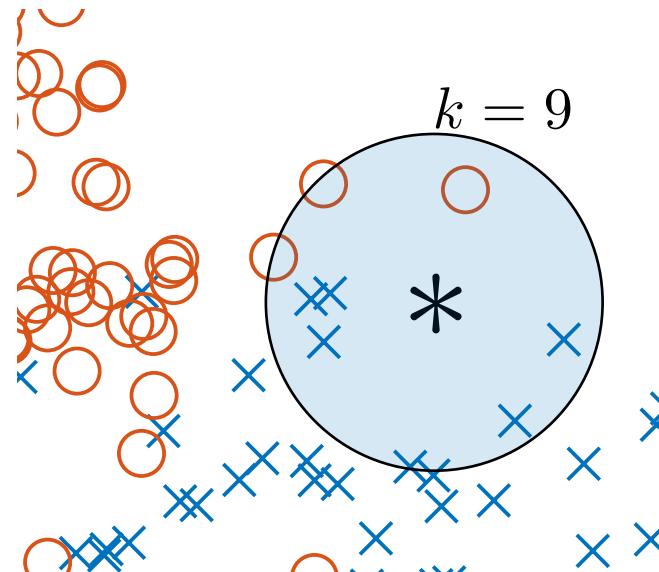
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

- k-Nearest Neighbors

- So, how do we know the two volumes of the kernels formed from different classes?
- The k NN classification algorithm
 1. For a test sample, find its k NN
 2. Count the number of training samples from each class
 - e.g. In the figure,
6 from class 1,
3 from class 2
 3. The class with more samples is denser
in that area → its kernel volume is smaller
→ choose the class with more samples

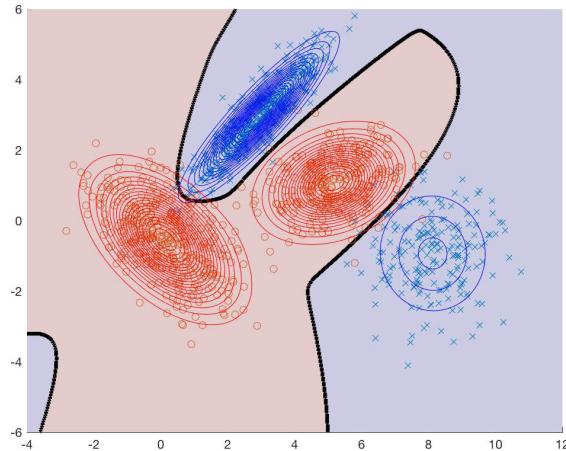


INDIANA UNIVERSITY

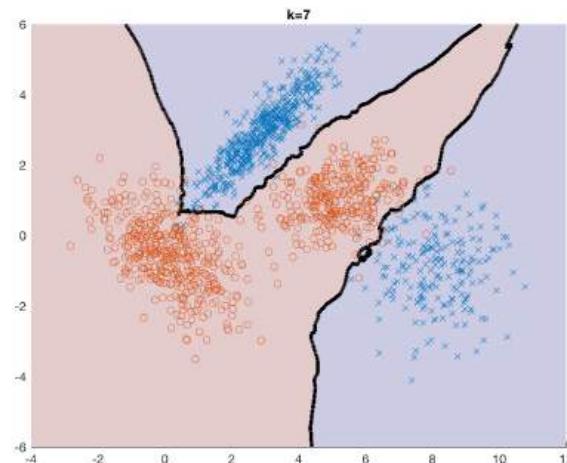
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Non-parametric Classifier

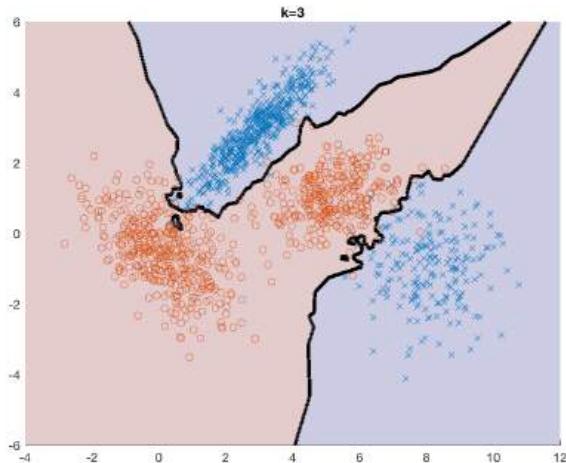
- k-Nearest Neighbors



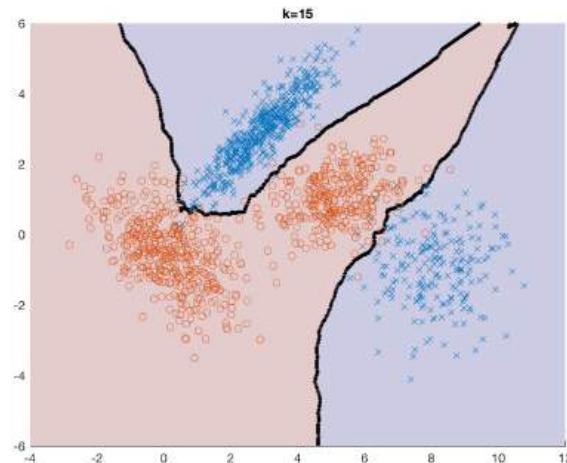
GMM Classifier



7NN Classifier



3NN Classifier



15NN Classifier



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Occam's Razor

- There is a trade-off between the model complexity and the performance of the classifier
 - Naïve Bayes classifier with isotropic Gaussians is easier to train than the one with full covariance matrix
 - But, if we succeed to estimate the full cov matrices, the decision boundaries are quite flexible
 - If you don't have enough training samples and fail to estimate the model parameters correctly, your decision boundaries will be messier
 - Are more complex classifiers always better than simpler ones?
 - Maybe not
 - Overfitting
- Non-parametric methods work well near the well sampled area
 - It's free from parameter estimation, but not entirely parameter free
 - e.g. the number of neighbors
 - Sorting is the bottleneck

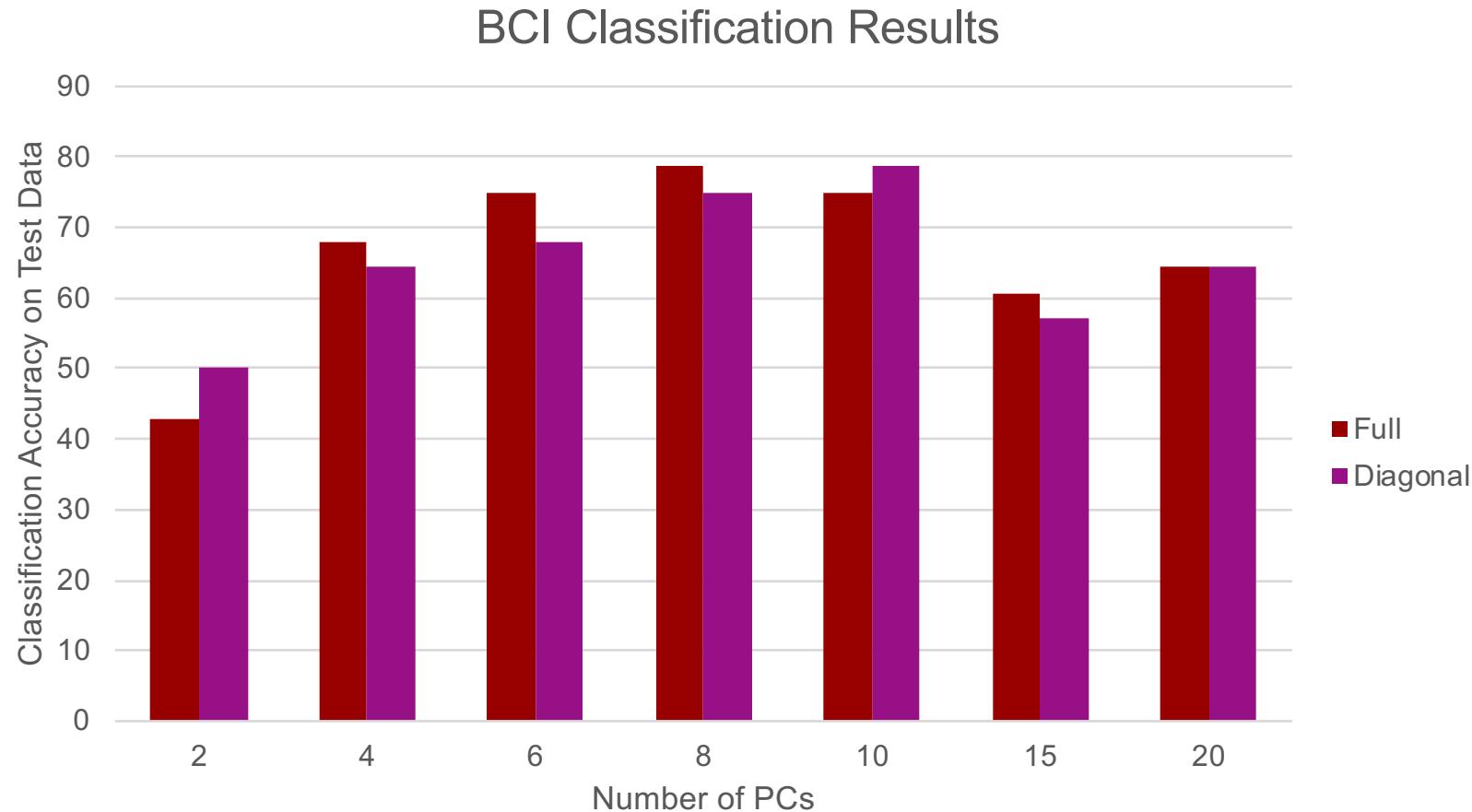


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

BCI Classification Results

- Left or right



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Reading

- Textbook 2.1 – 2.6



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



Thank You!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING