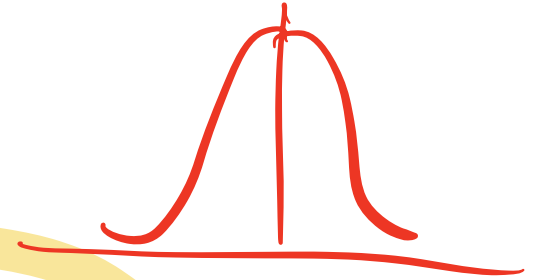We will use the Bayesian decision criteria applied to normally distributed classes, whose parameters are either known or estimated from the sample.

# Parametric Classification

# Parametric Classification

- If $p(\mathbf{x} \mid C_i) = \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

$$p(\mathbf{x} \mid C_i) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_i|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right]$$
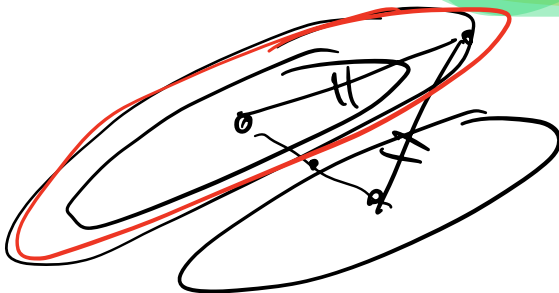
- Discriminant functions are:

*→ think as score for $C_i$.*

$$g_i(\mathbf{x}) = P(C_i \mid \mathbf{x})$$

*$\downarrow$ i=1 to K classes*

$$g_i(\mathbf{x}) = \log P(C_i \mid \mathbf{x}) \quad = \log\left(P(\mathbf{x} \mid C_i) \cdot P(C_i)\right)$$

$$= \log p(\mathbf{x} \mid C_i) + \log P(C_i)$$

$$= -\frac{d}{2}\log 2\pi - \frac{1}{2}\log|\boldsymbol{\Sigma}_i| - \frac{1}{2}(\mathbf{x} - \mu_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \mu_i) + \log P(C_i)$$
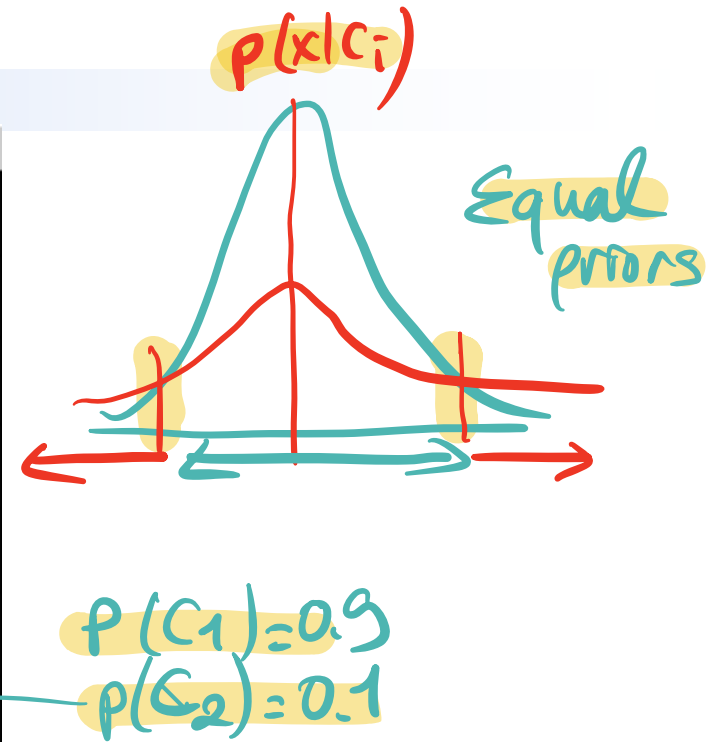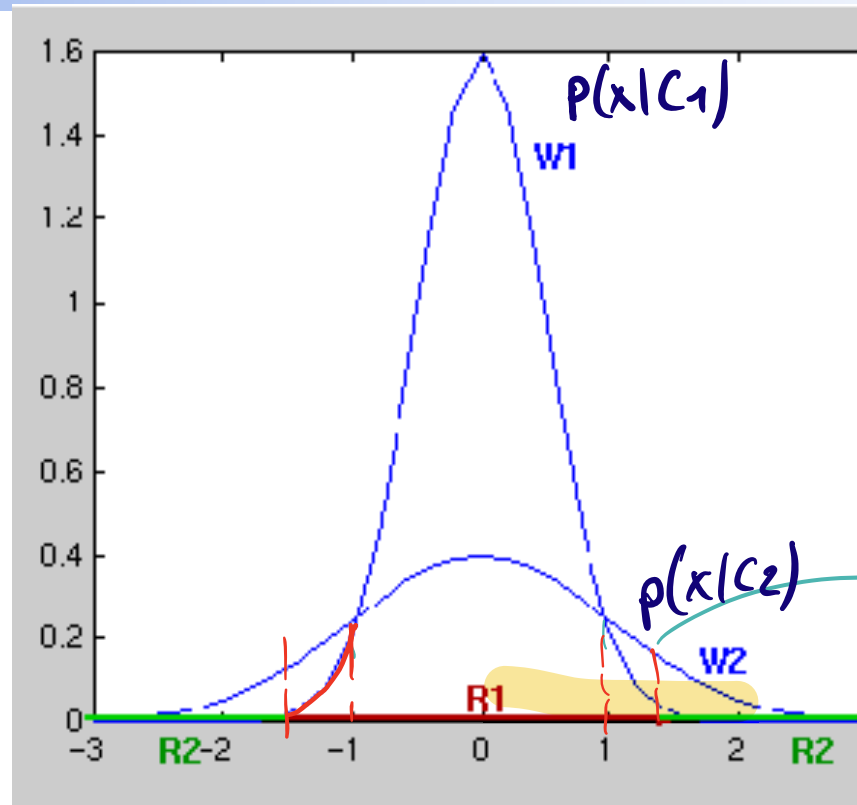
2

# Estimation of Parameters

If we estimate the unknown parameters from the sample, the discriminant function becomes:

$$g_i(\mathbf{x}) = -\frac{1}{2}\log|\mathbf{S}_i| - \frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

$\mu, \sigma$: population parameters

$M, S$: sample $''$

*Typical single-variable normal distributions showing a disconnected decision region $R_2$*

Handwritten annotations:

$p(x|C_i)$

Equal priors

$p(C_1) = 0.9$
$p(C_2) = 0.1$

$$\exp\left\{\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

In figure: $p(x|C_1)$, W1, $p(x|C_2)$, W2, R1, R2

- To illustrate the previous result, we will compute the decision boundaries for a 3-class, 2-dimensional problem with the following class mean vectors and covariance matrices and equal priors
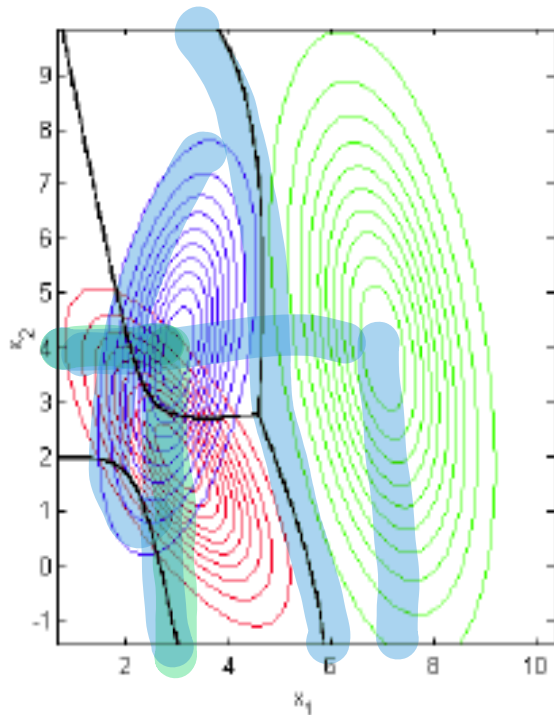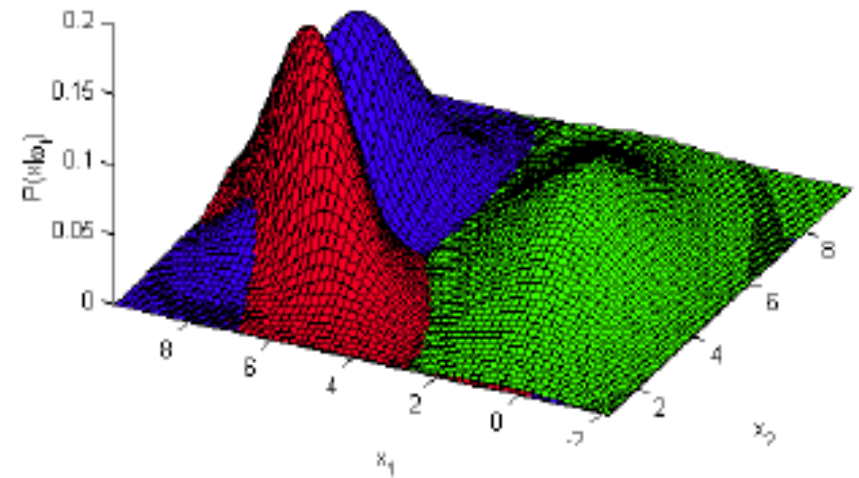
$$\mu_1 = \begin{bmatrix} 3 & 2 \end{bmatrix}^T \qquad \mu_2 = \begin{bmatrix} 5 & 4 \end{bmatrix}^T \qquad \mu_3 = \begin{bmatrix} 2 & 5 \end{bmatrix}^T$$

$$\Sigma_1 = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \qquad \Sigma_2 = \begin{bmatrix} 1 & -1 \\ -1 & 7 \end{bmatrix} \qquad \Sigma_3 = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 3 \end{bmatrix}$$
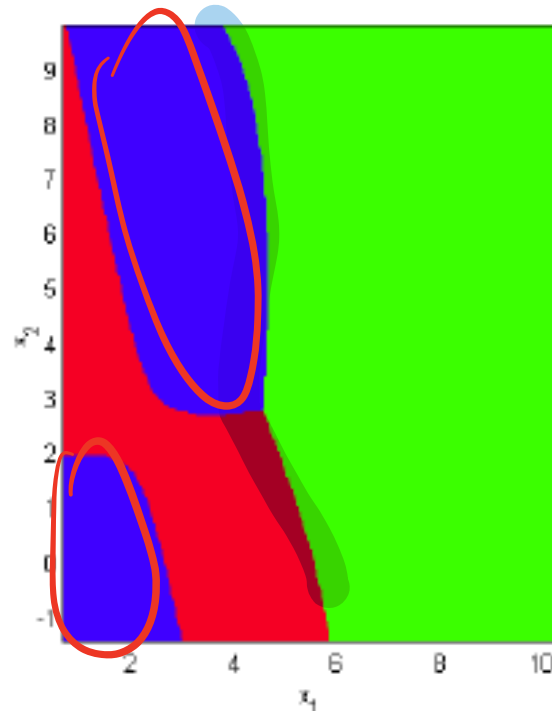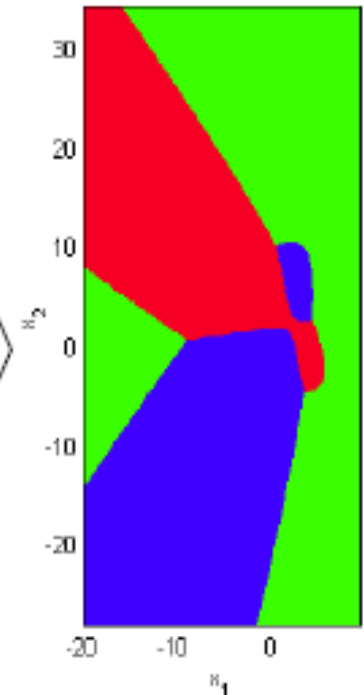
red          green          blue



Boundaries

Zoom out

5

- If d (dimension) is large with respect to N (number of samples), we may have a problem with this approach:
  - $|\Sigma|$ may be zero, thus $\Sigma$ will be singular (inverse does not exist)
  - $|\Sigma|$ may be non-zero, but very small, instability would result
    - Small changes in $\Sigma$ would cause large changes in $\Sigma^{-1}$

- Solutions:
  - Reduce the dimensionality
    - Feature selection
    - Feature extraction: PCA,…
  - Pool the data and estimate a common covariance matrix for all classes

$$\Sigma = \Sigma_i\, P(C_i) * \Sigma_i$$

$\Sigma$

weight

$$\sigma_{12} = \frac{1}{N} \sum_{i=1}^{N} (x_1^{(i)} - \mu_1)(x_2^{(i)} - \mu_2)$$

weight

- In the following slides, **we will make increasing assumptions about the covariance matrix** and see what the corresponding discriminant function and resulting boundaries look like.
  - QDA, LDA, Naïve Bayes, Nearest Mean classifiers

# Case 2) Common Covariance Matrix $S=S_i$

$$S_i = S$$

- **Shared common sample covariance S**
  - □ An arbitrary covariance matrix – **but shared between the classes**

- We had this full discriminant function:

$$g_i(\mathbf{x}) = -\frac{1}{2}\log|\mathbf{S}| - \frac{1}{2}(\mathbf{x}-\mathbf{m}_i)^T \mathbf{S}_i^{-1}(\mathbf{x}-\mathbf{m}_i) + \log \hat{P}(C_i)$$

which now reduces to:

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x}-\mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x}-\mathbf{m}_i) + \log \hat{P}(C_i)$$

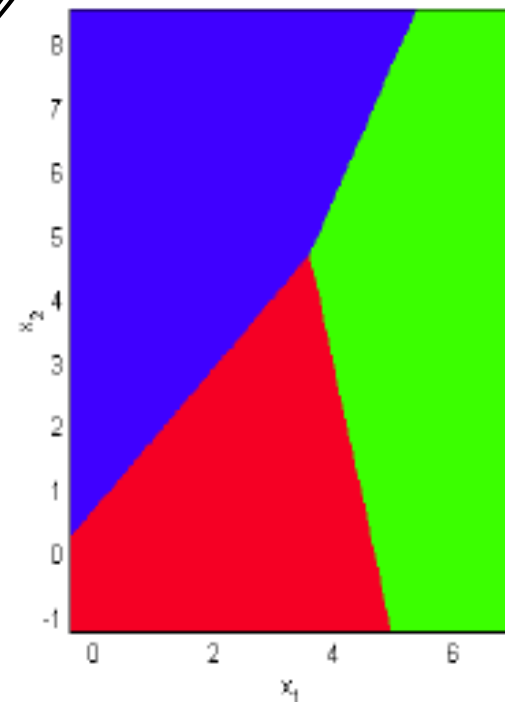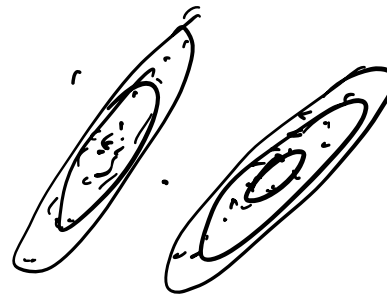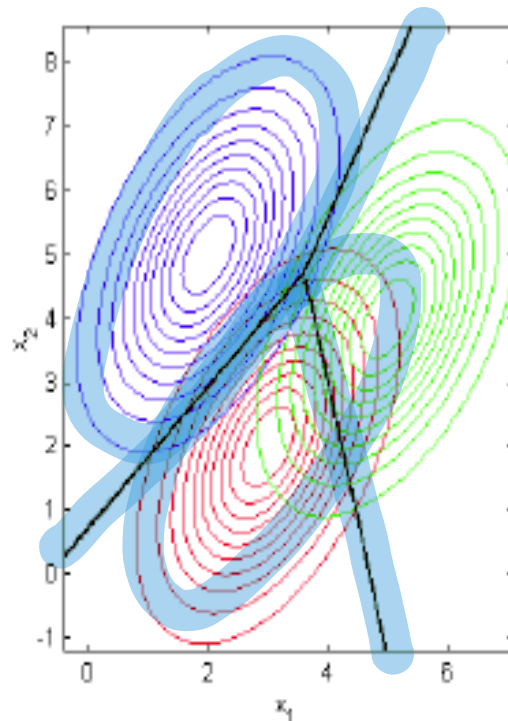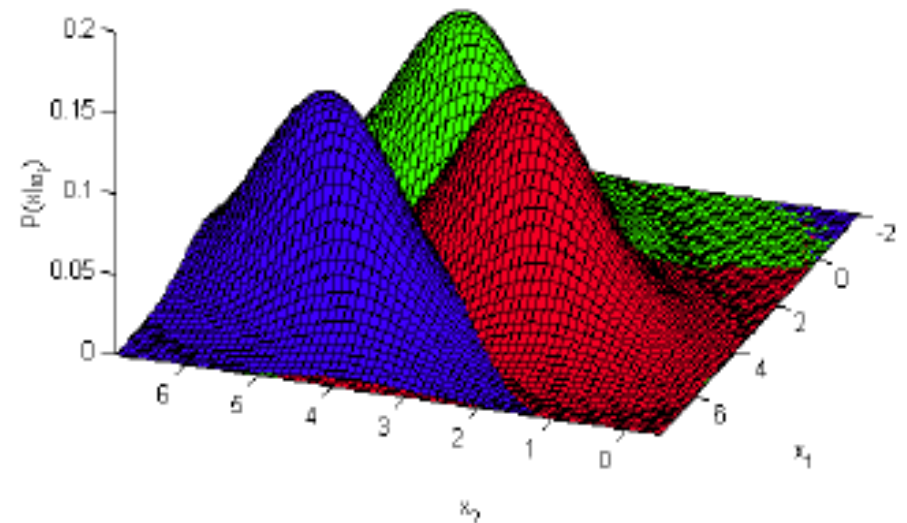which is a **linear discriminant** (decision boundaries are hyper-planes)

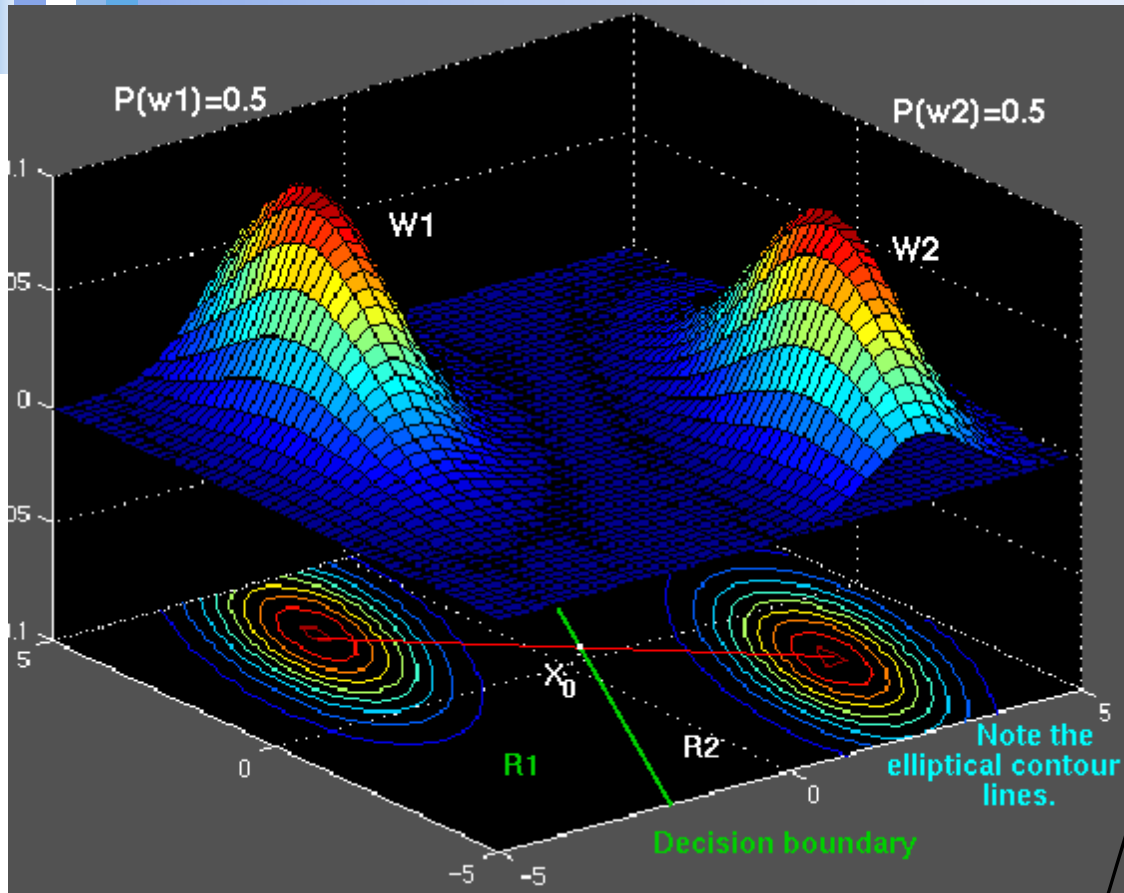Which class to assign $x$ to ? $\max_i g_i(x)$

9

- To illustrate the previous result, we will compute the decision boundaries for a 3-class, 2-dimensional problem with the following class mean vectors and covariance matrices and equal priors
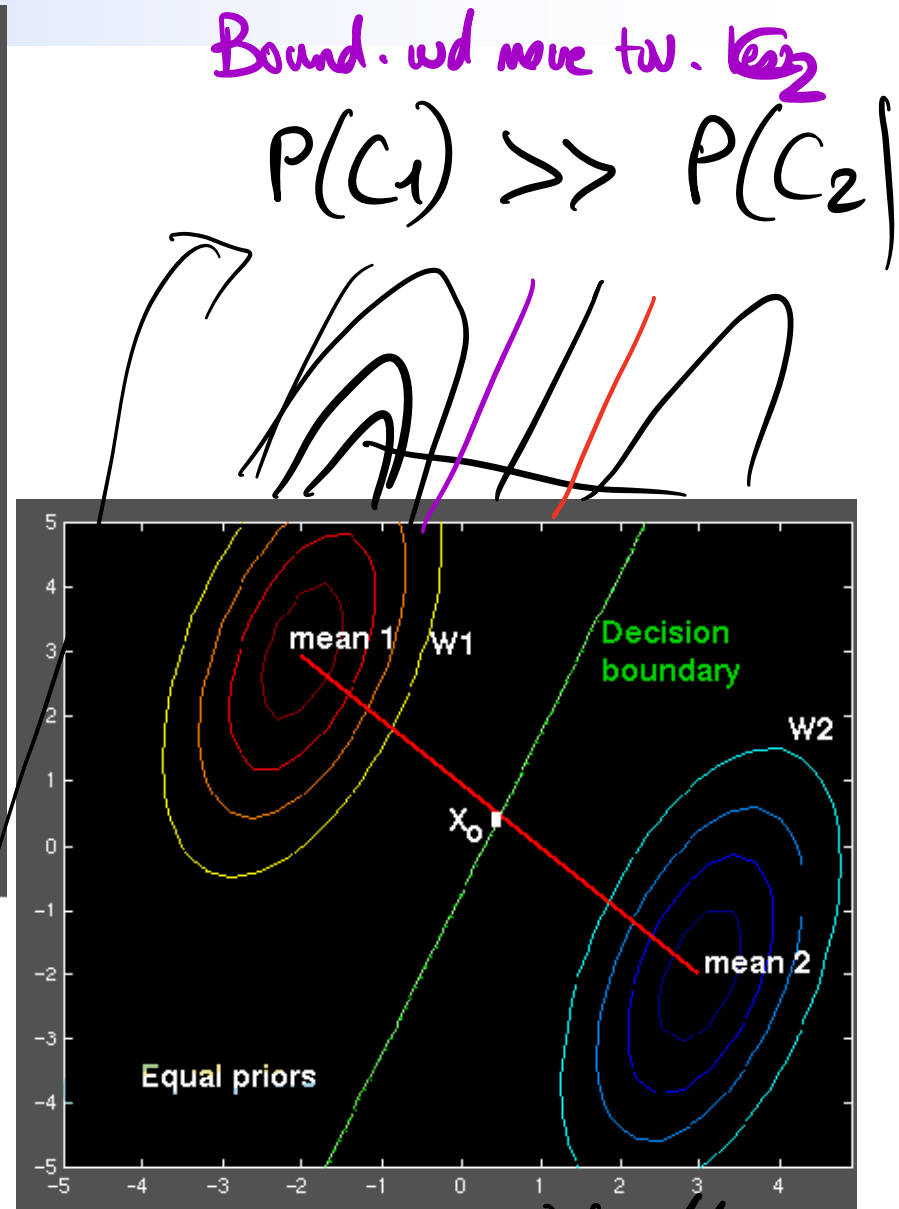
$$\mu_1 = \begin{bmatrix} 3 & 2 \end{bmatrix}^T \quad \mu_2 = \begin{bmatrix} 5 & 4 \end{bmatrix}^T \quad \mu_3 = \begin{bmatrix} 2 & 5 \end{bmatrix}^T$$

$$\Sigma_1 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix} \quad \Sigma_3 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 2 \end{bmatrix}$$
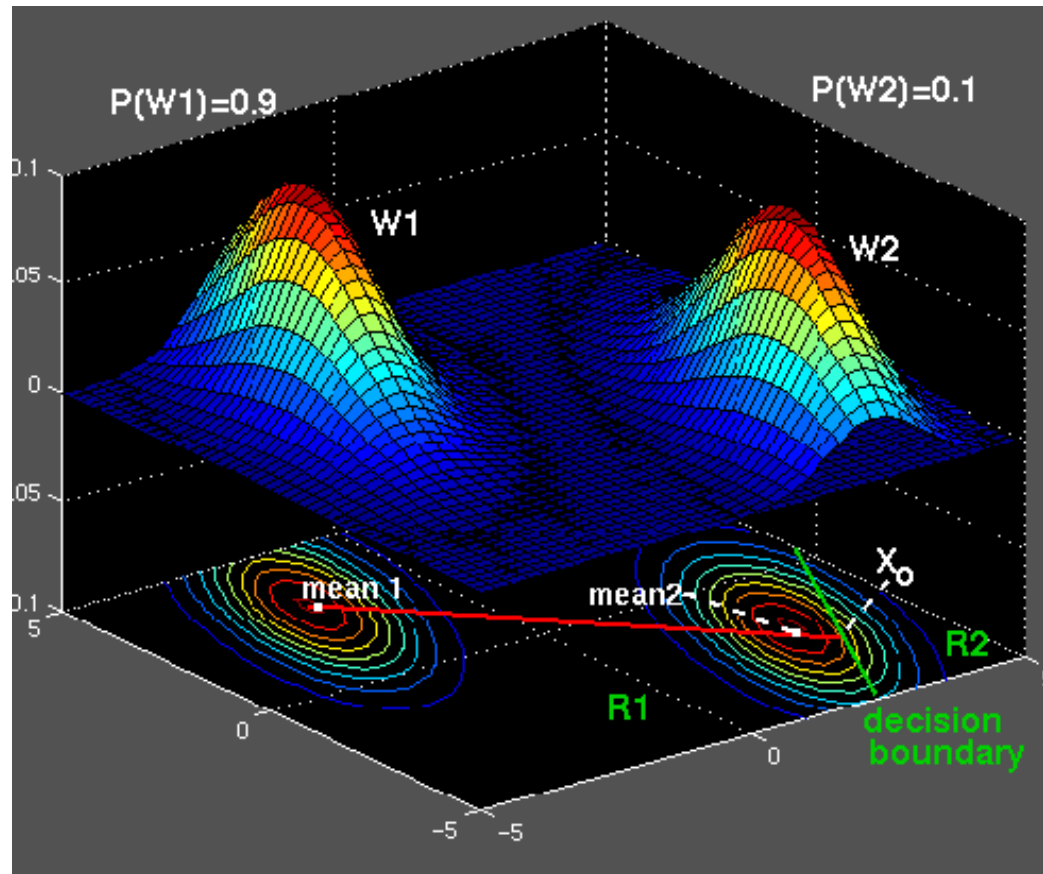
**If we assume equal class priors**, the classifier becomes a minimum Mahalanobis classifier

Bound. wd move tw. less₂

$$P(C_1) \gg P(C_2)$$

Not equal priors! what happens -

13

Unequal priors shift the decision boundary towards the less likely class.

# Case 3) Common Covariance Matrix *S* which is *Diagonal*

$$S_i = S + Diagonal$$

- In the previous case, we had a common, general covariance matrix, resulting in these discriminant functions:

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

$$S = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \phi \\ \phi & & \sigma_3^2 \end{bmatrix}$$

- **When $x_j$ ($j = 1,..d$) are independent** (or assumed to be independent for simplicity), **then $\Sigma$ is diagonal:**

*15*

# *Case 3) Common Covariance Matrix $S$ which is Diagonal*

- In the previous case, we had a common, general covariance matrix, resulting in these discriminant functions:

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$
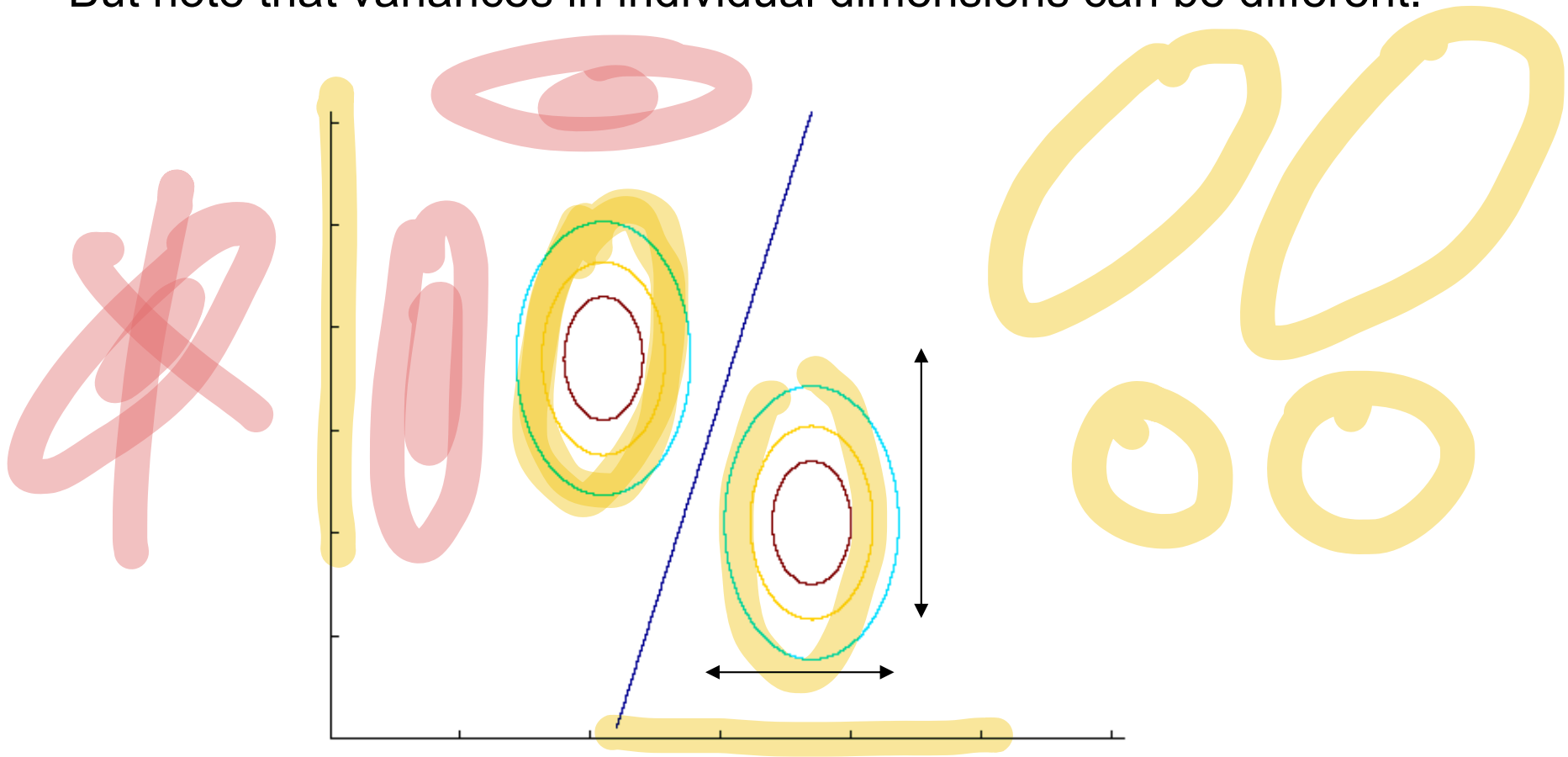
- When $x_j$ $(j = 1,..d)$ are independent (or assumed to be independent for simplicity), then $\Sigma$ is diagonal.

This is the **Naive Bayes classifier** where $p(x_j|C_i)$ are univariate Gaussian.
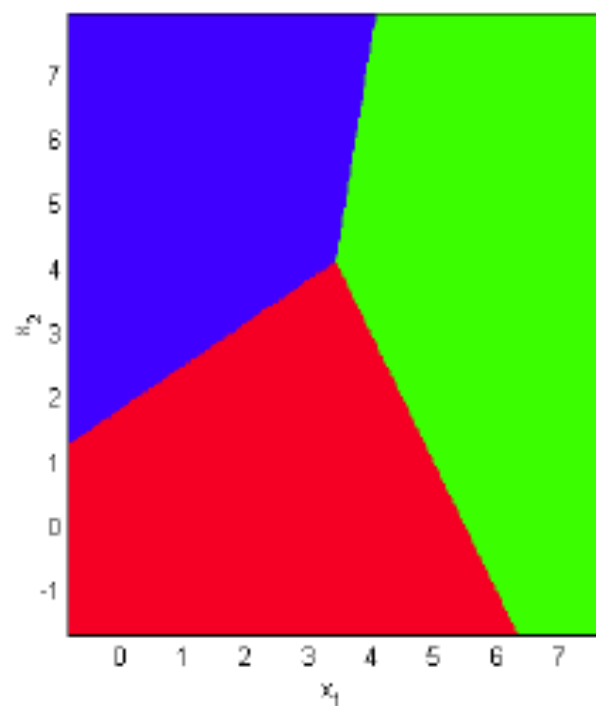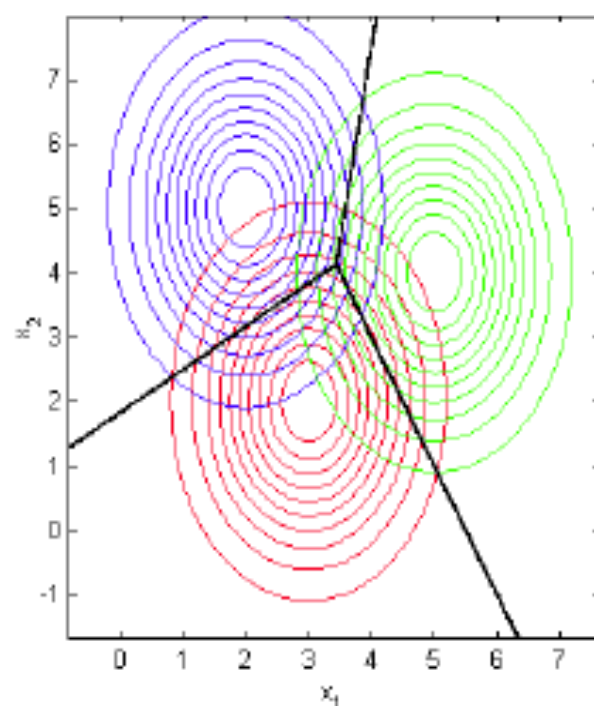
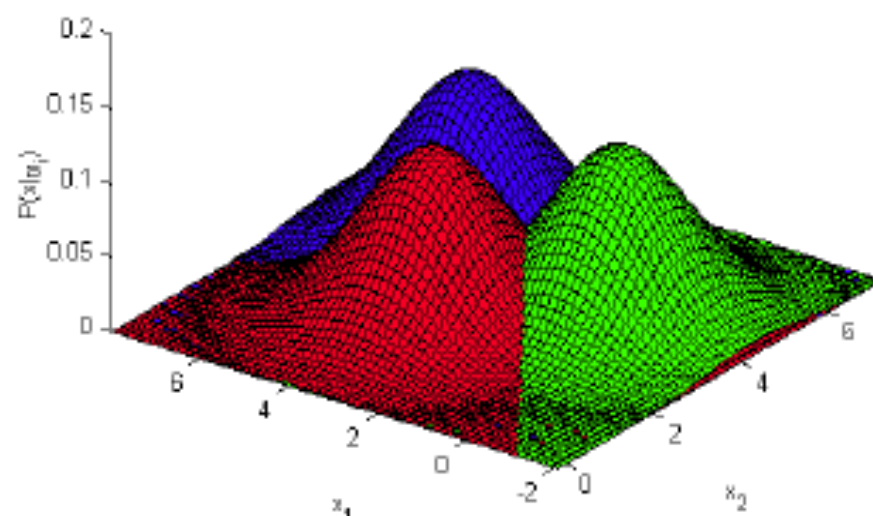# Case 3) Common Covariance Matrix $S$ which is Diagonal

Diagonal covariance matrices means no correlation among attributes – hence contours are axis-aligned

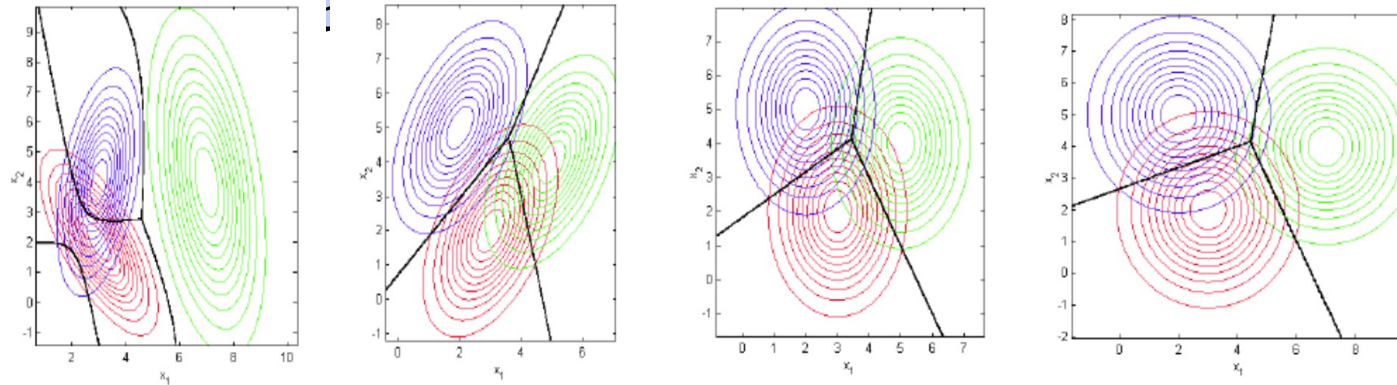But note that variances in individual dimensions can be different.

- To illustrate the previous result, we will compute the decision boundaries for a 3-class, 2-dimensional problem with the following class mean vectors and covariance matrices and equal priors

$$\mu_1 = \begin{bmatrix} 3 & 2 \end{bmatrix}^T \quad \mu_2 = \begin{bmatrix} 5 & 4 \end{bmatrix}^T \quad \mu_3 = \begin{bmatrix} 2 & 5 \end{bmatrix}^T$$

$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad \Sigma_3 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$
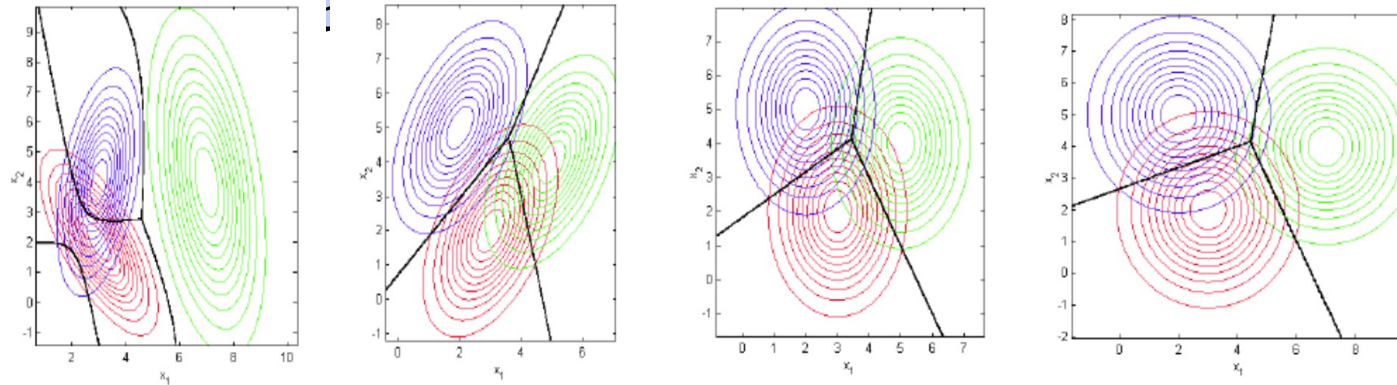
# Model Selection



| Assumption | Covariance matrix | Number of parameters |
|---|---|---|
| Case 1) None. All different, Hyperellipsoidal | $\mathbf{S}_i$ | $K\,d(d+1)/2$ |
| Case 2) Shared, Hyperellipsoidal | $\mathbf{S}_i=\mathbf{S}$ | $\cancel{K}\;d(d+1)/2$ |
| Case 3) Shared & Axis-aligned | $\mathbf{S}_i=\mathbf{S}$, with $s_{ij}=0$ | $d$ |
| Case 4) Shared & Hyperspheric | $\mathbf{S}_i=\mathbf{S}=s^2\mathbf{I}$ | 1 |

- As we increase complexity (less restricted **S**), it is more important to have sufficient data to properly estimate the parameters.
- Simpler models may not model the underlying distributions fully correctly, but may even work better.
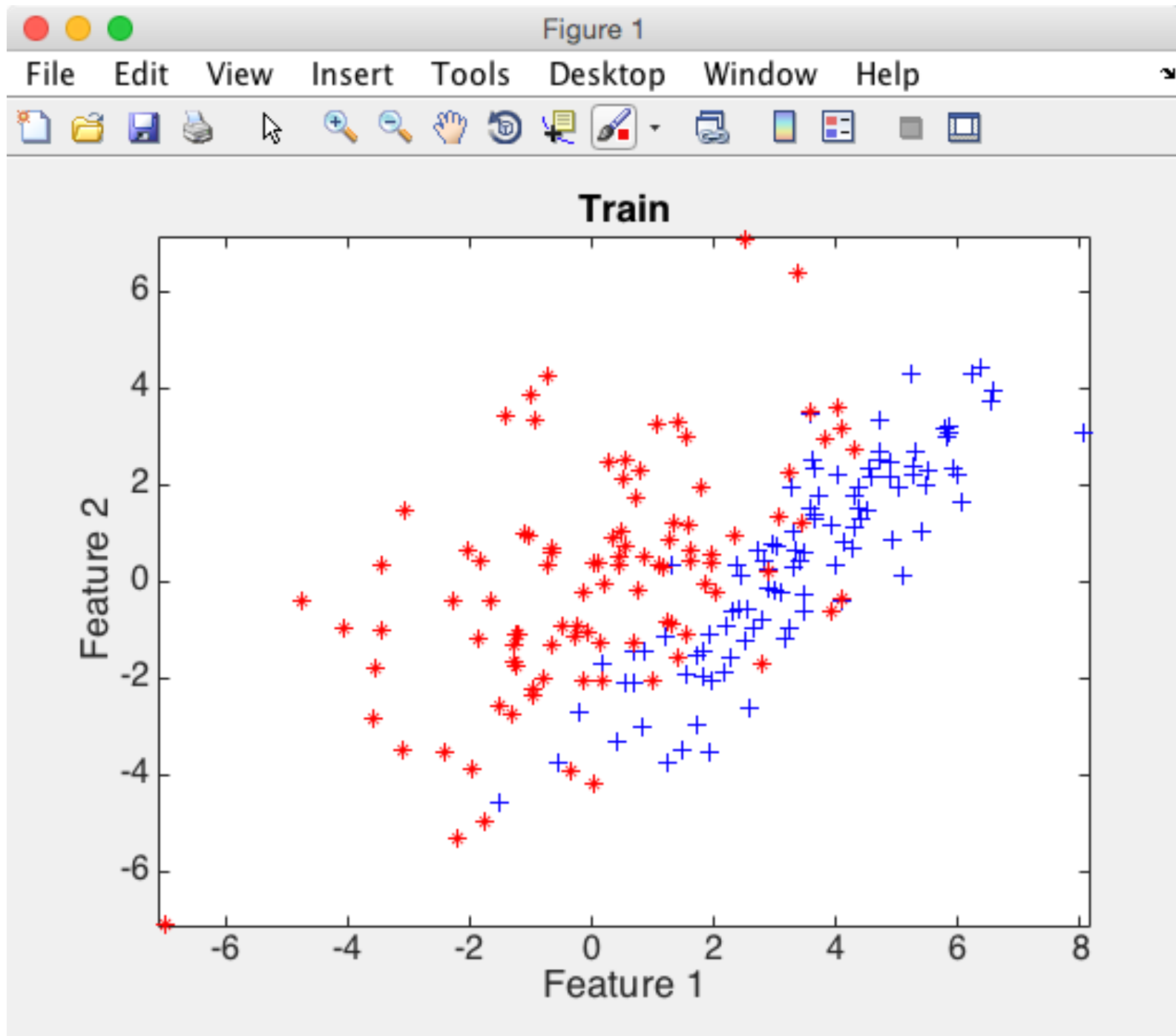
19

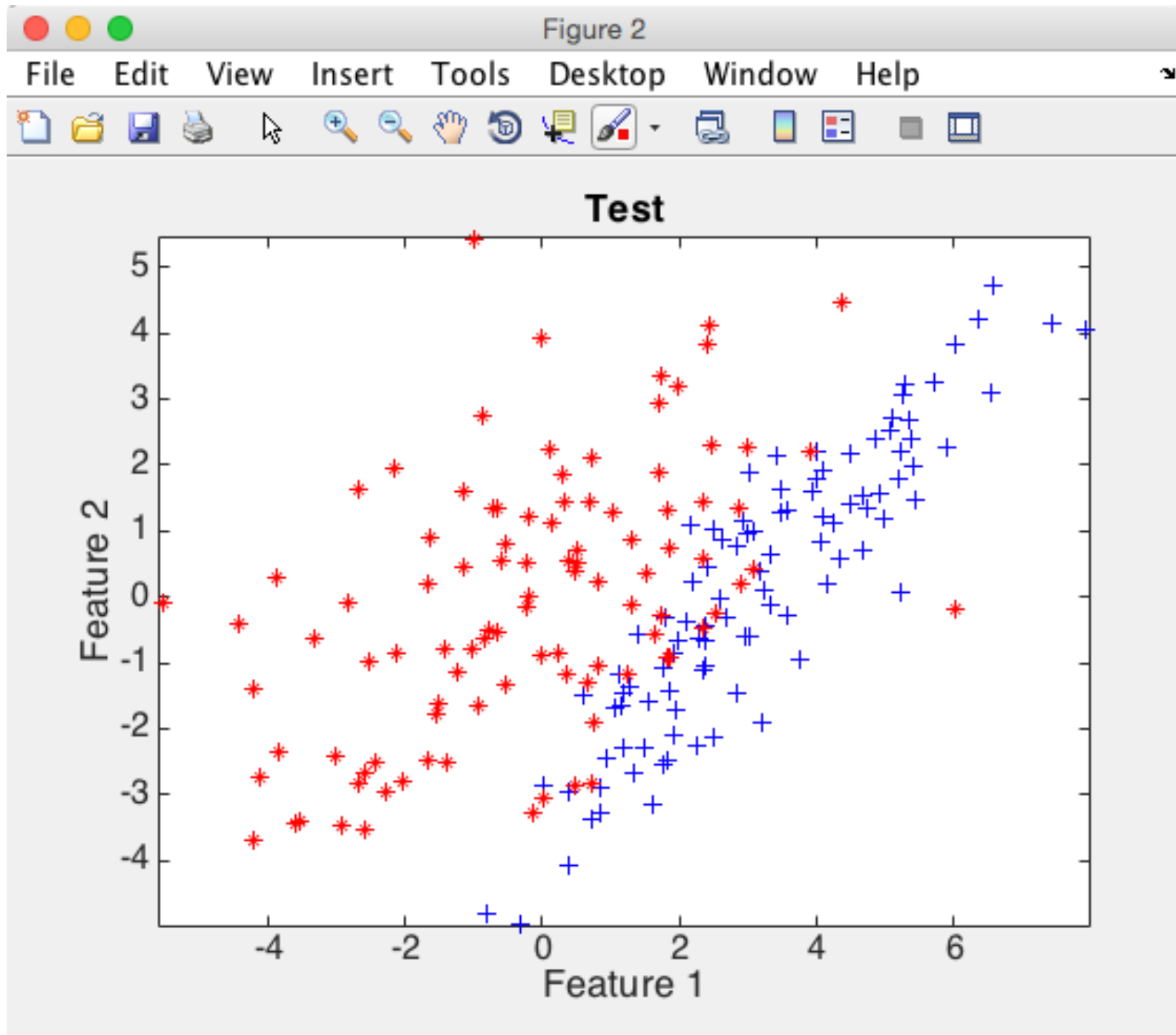| Assumption | Covariance matrix | Number of parameters |
|---|---|---|
| Case 1) None. All different, Hyperellipsoidal | $\mathbf{S}_i$ | $K\,d(d+1)/2$ |
| Case 2) Shared, Hyperellipsoidal | $\mathbf{S}_i = \mathbf{S}$ | $d(d+1)/2$ |
| Case 3) Shared & Axis-aligned | $\mathbf{S}_i = \mathbf{S}$, with $s_{ij}=0$ | $d$ |
| Case 4) Shared & Hyperspheric | $\mathbf{S}_i = \mathbf{S} = s^2\mathbf{I}$ | $1$ |

- As we increase complexity (less restricted **S**), bias decreases and variance increases
- Assume simple models (allow some bias) to control variance (regularization)

- **QDA** (short for Quadratic Bayes classifier or Quadratic Discriminant Analysis) and

- **LDA** (short for Linear Bayes classifier or Linear Discriminant Analysis) are the two Gaussian Bayes classifiers;….

  - First one corresponds to the general covariance matrix case and second one to the shared covariance matrix case

  - See http://scikit-learn.org/stable/modules/lda_qda.html Section 1.2.2.
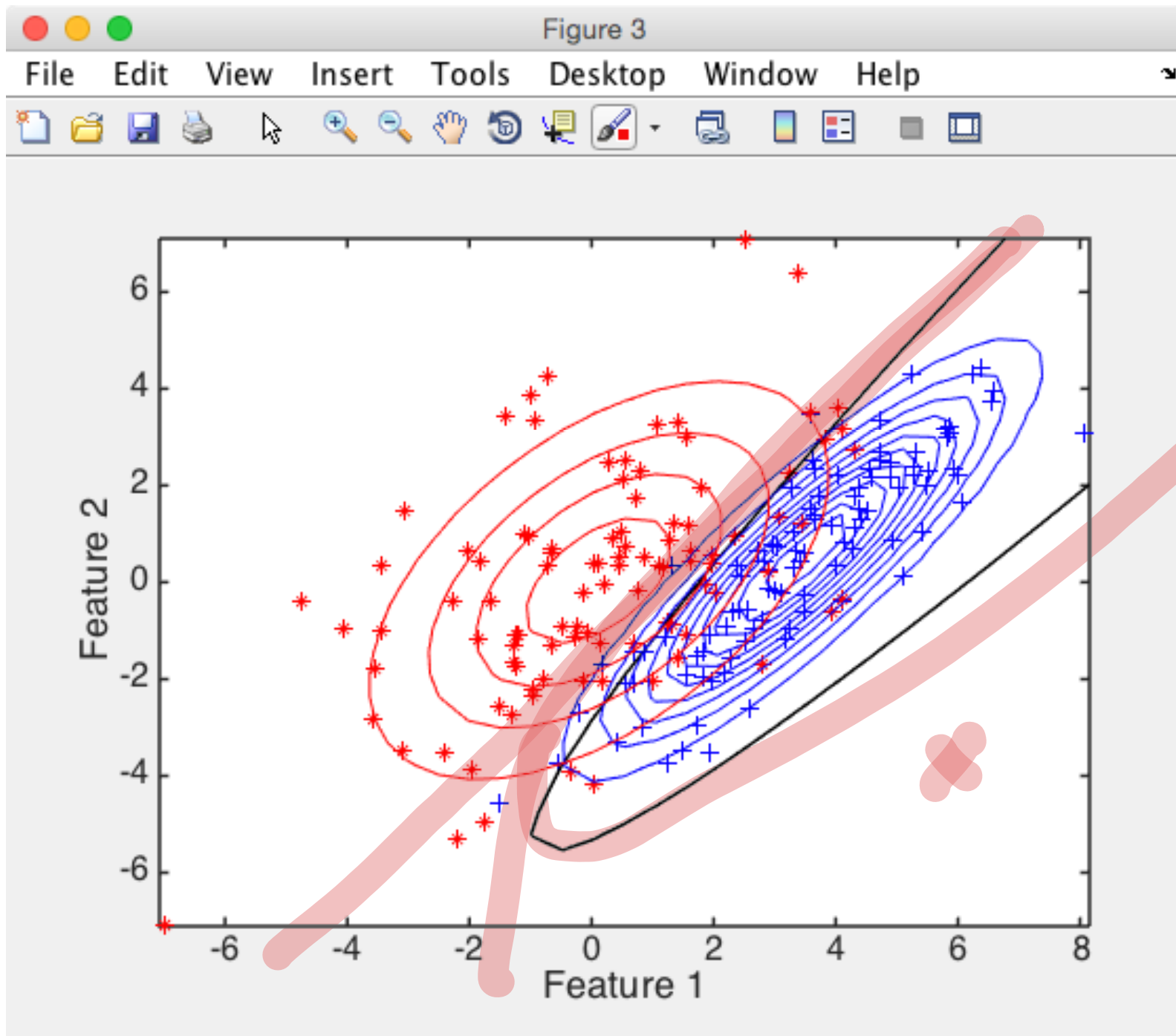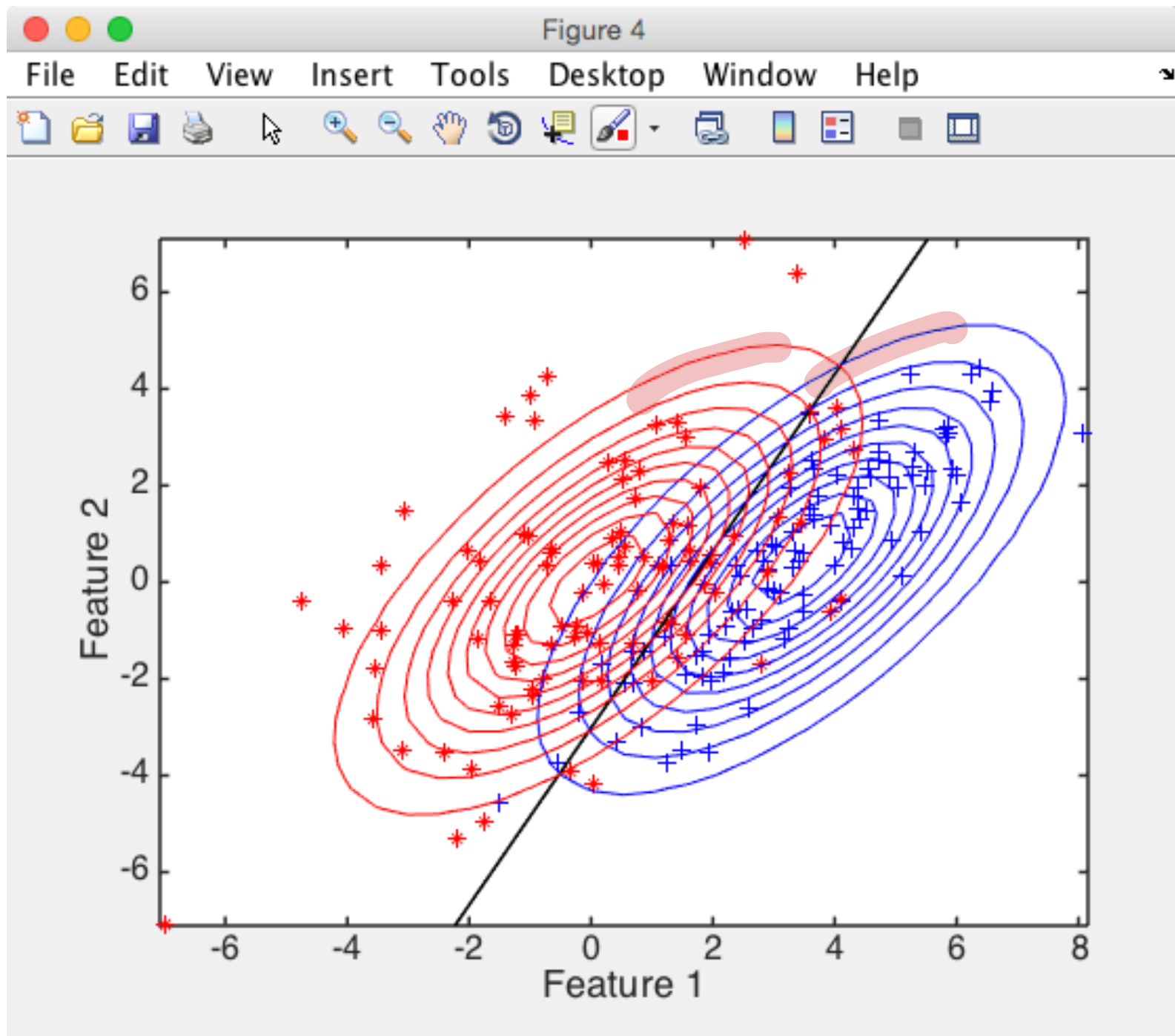
# *Matlab Exercise: Generate some training data*

# *Matlab Exercise: Generate some test data*

# *Quadratic classifier*

32

# *Another example* with *N=20 points*