

# CHAPTER 5: MULTIVARIATE METHODS

**Stella Grasshof**

# Overview of today

- 1) reminder and hint
- 2) multivariate data
  - ▣ descriptive statistics
  - ▣ classification
  - ▣ model selection
- 3) regression
  - ▣ 1D regression revisited
  - ▣ multivariate regression
  - ▣ model selection

# Reminder and Hint

□ Pre-warning: many equations today (will be less Th.)

□ List of errors in the book (incomplete)

<https://www.cmpe.boun.edu.tr/~ethem/i2ml3e/>

□ Helpful online tool:

<https://www.wolframalpha.com/>

Solve an equation with parameters:

solve  $a x^2 + b x + c = 0$  for  $x$

=

More examples

More examples

## Simplification

Simplify algebraic functions and expressions.

Simplify an expression:

$1/(1+\sqrt{2})$

=

simplify  $x^5 - 20x^4 + 163x^3 - 676x^2 + 1424x - 1209$

=

simplify  $\cos(\arcsin(x)/2)$

=

## Rational Functions

Compute discontinuities and other properties of rational functions.

Compute properties of a rational function:

$(x^2 - 1)/(x^2 + 1)$

=

Compute a partial fraction decomposition:

partial fractions  $(x^2 - 4)/(x^4 - x)$

=

More examples

- Arithmetic
- Calculus & Analysis
- Geometry
- Linear Algebra

## Matrices

Find properties and perform computations on matrices.

Do basic arithmetic on matrices:

$\{0, -1\}, \{1, 0\} \cdot \{1, 2\}, \{3, 4\} + \{2, -1\}, \{-1, 2\}$

=

Compute eigenvalues and eigenvectors of a matrix:

eigenvalues  $\{4, 1\}, \{2, -1\}$

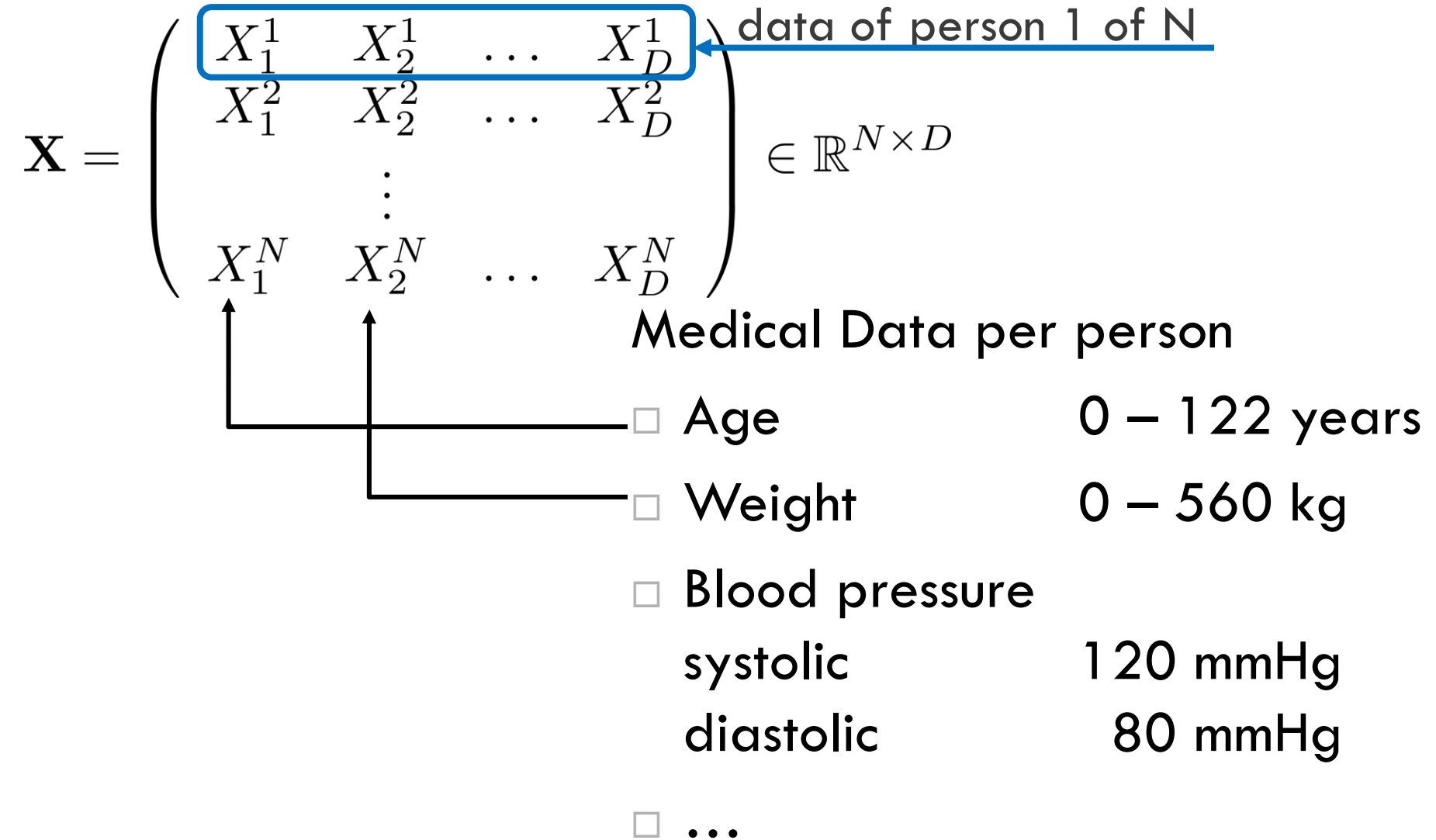
=

# Multivariate Data

- Multiple measurements (sensors)
- $D$  inputs/features/attributes:  $D$ -variate
- $N$  instances/observations/examples

$$\mathbf{X} = \begin{pmatrix} X_1^1 & X_2^1 & \dots & X_D^1 \\ X_1^2 & X_2^2 & \dots & X_D^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^N & X_2^N & \dots & X_D^N \end{pmatrix} \in \mathbb{R}^{N \times D}$$

# Multivariate Data



# Multivariate Parameters

$$\mathbf{X} = \begin{pmatrix} X_1^1 & X_2^1 & \dots & X_D^1 \\ X_1^2 & X_2^2 & \dots & X_D^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^N & X_2^N & \dots & X_D^N \end{pmatrix} \in \mathbb{R}^{N \times D}$$

□ **Mean**

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_D]^T \in \mathbb{R}^D$$

□ **Covariance**  $\sigma_{ij} = \text{Cov}(X_i, X_j)$   $X_i \in \mathbb{R}^N$  column  $i$  of matrix  $\mathbf{X}$

$$\boldsymbol{\Sigma} \equiv \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1D} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{D1} & \sigma_{D2} & \dots & \sigma_D^2 \end{pmatrix} \in \mathbb{R}^{D \times D}$$

$$\boldsymbol{\Sigma} \equiv \text{Cov}(\mathbf{X}) = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$$

□ **Correlation**  $\text{Corr}(X_i, X_j) = \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$

**Matrix minus vector**

# Multivariate Parameters

## Adapted Notation

vector of random variables

$$\mathbf{x} = (x_1, \dots, x_D)^T \in \mathbb{R}^D$$

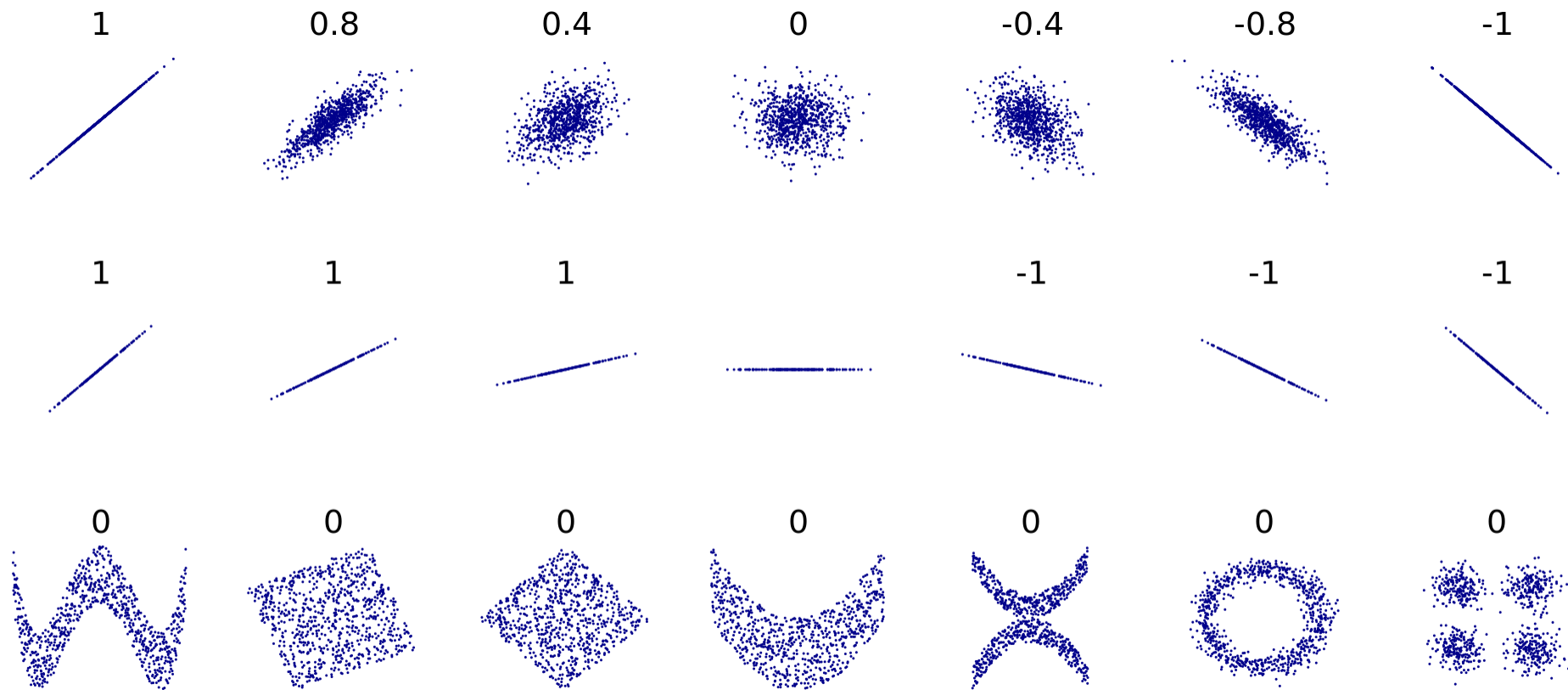
□ **Mean**  $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_D]^T \in \mathbb{R}^D$

□ **Covariance**  $\sigma_{ij} = \text{Cov}(x_i, x_j)$

$$\boldsymbol{\Sigma} \equiv \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1D} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{D1} & \sigma_{D2} & \dots & \sigma_D^2 \end{pmatrix} \in \mathbb{R}^{D \times D}$$
$$\boldsymbol{\Sigma} \equiv \text{Cov}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]$$

□ **Correlation**  $\text{Corr}(x_i, x_j) = \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j} \in [-1, 1]$

# Correlation



correlation  $\neq$  causality



# Correlation $\neq$ Causality

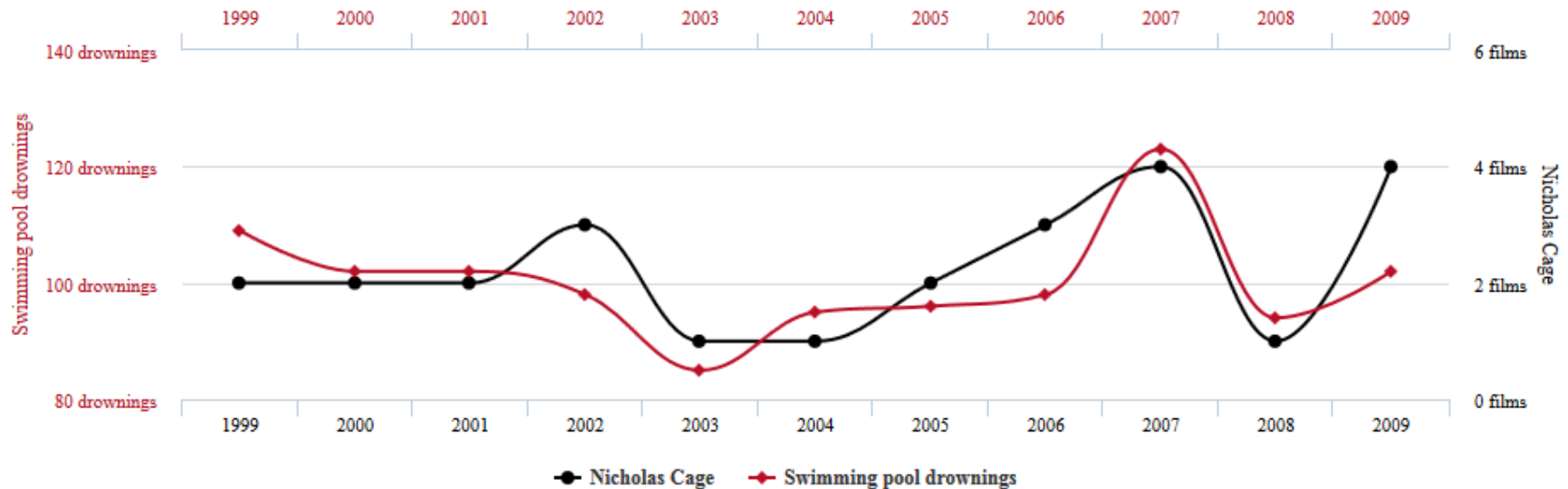
see: <https://www.tylervigen.com/spurious-correlations>

**Number of people who drowned by falling into a pool**

correlates with

**Films Nicolas Cage appeared in**

Correlation: 66.6% ( $r=0.666004$ )



Data sources: Centers for Disease Control & Prevention and Internet Movie Database

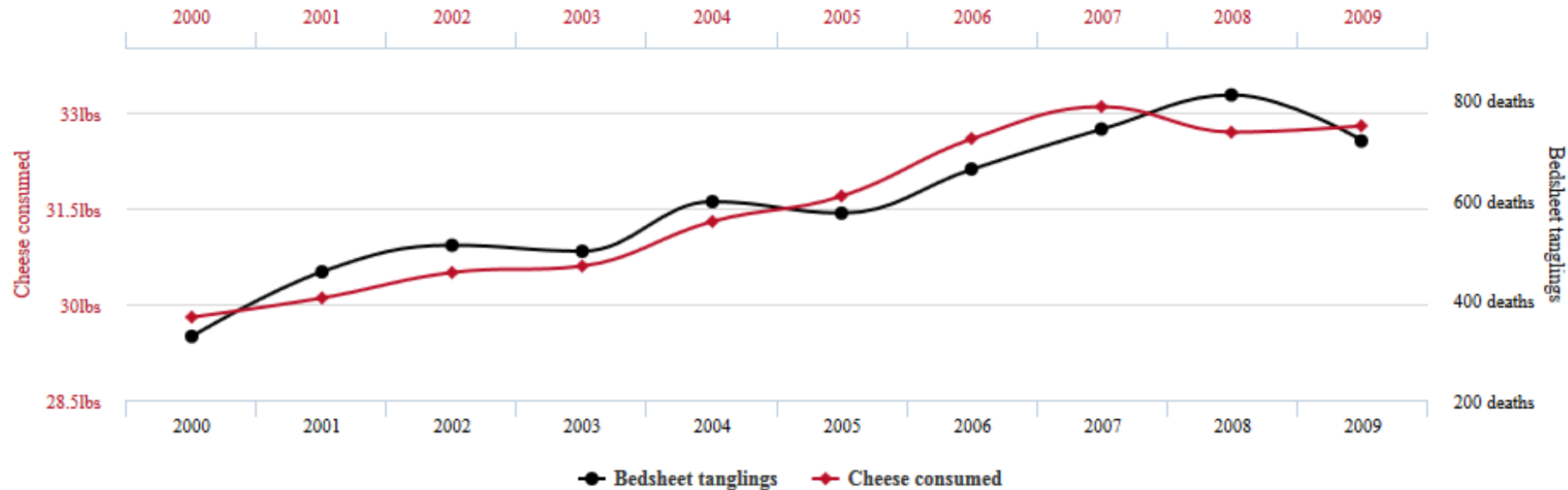
tylervigen.com

# Correlation $\neq$ Causality

see: <https://www.tylervigen.com/spurious-correlations>

**Per capita cheese consumption**  
correlates with  
**Number of people who died by becoming tangled in their bedsheets**

Correlation: 94.71% ( $r=0.947091$ )



# Uncorrelated $\neq$ Independence

- $x, y$  uncorrelated, i. e.  $\text{Corr}(x, y) = \rho_{xy} = 0$   
if they are not **linearly** related
- $x, y$  are independent if joint p.d.f. can be  
factorized as  $p(x, y) = p_x(x)p_y(y)$

$x, y$  independent  $\Rightarrow$   $x, y$  uncorrelated  
 $\nLeftarrow$

$\Leftarrow$   
**ONLY IF  $x, y$  are normally distributed**

$$x \sim \mathcal{N}(\mu_x, \sigma_x^2), y \sim \mathcal{N}(\mu_y, \sigma_y^2)$$

# Parameter Estimation from Data

$$\mathbf{X} = \begin{pmatrix} X_1^1 & X_2^1 & \dots & X_D^1 \\ X_1^2 & X_2^2 & \dots & X_D^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^N & X_2^N & \dots & X_D^N \end{pmatrix} \in \mathbb{R}^{N \times D}$$

□ **Sample Mean**  $\mu \approx \hat{\mu} = \mathbf{m} = \frac{1}{N} \left( \sum_{n=1}^N X_1^n, \dots, \sum_{n=1}^N X_D^n \right)^T \in \mathbb{R}^D$

□ **Covariance**  $\sigma_{ij} \approx \hat{\sigma}_{ij} = s_{ij} = \frac{1}{N} \sum_{n=1}^N (X_i^n - m_i)(X_j^n - m_j)$

$$s_i^2 = \frac{1}{N} \sum_{n=1}^N (X_i^n - m_i)^2$$

□ **Correlation**  $\rho_{ij} \approx \hat{\rho}_{ij} = r_{ij} = \frac{s_{ij}}{s_i s_j}$

# Parameter Estimation from Data

$$\mathbf{X} = \begin{pmatrix} X_1^1 & X_2^1 & \dots & X_D^1 \\ X_1^2 & X_2^2 & \dots & X_D^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^N & X_2^N & \dots & X_D^N \end{pmatrix} \in \mathbb{R}^{N \times D}$$

□ **Sample Mean**  $\mu \approx \hat{\mu} = \mathbf{m} = \frac{1}{N} \left( \sum_{n=1}^N X_1^n, \dots, \sum_{n=1}^N X_D^n \right)^T \in \mathbb{R}^D$

□ **Covariance**  $\sigma_{ij} \approx \hat{\sigma}_{ij} = s_{ij} = \frac{1}{N-1} \sum_{n=1}^N (X_i^n - m_i)(X_j^n - m_j)$

unbiased estimator

$$s_i^2 = \frac{1}{N-1} \sum_{n=1}^N (X_i^n - m_i)^2$$

□ **Correlation**  $\rho_{ij} \approx \hat{\rho}_{ij} = r_{ij} = \frac{s_{ij}}{s_i s_j}$

# Missing Values ?

What to do if certain instances have missing attributes?

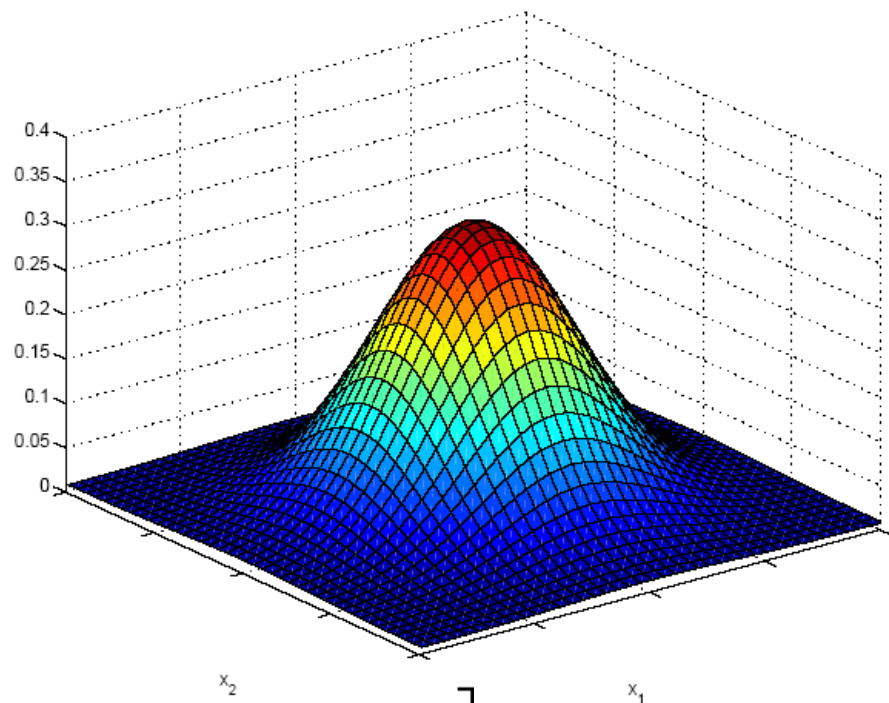
- **Ignore** those instances:  
not a good idea if the sample is small
- **Use 'missing' as an attribute:**  
may give information
- **Imputation:** Fill in the missing value
  - ▣ Mean imputation: Use the most likely value (e.g., mean)
  - ▣ Imputation by regression:  
Predict based on other attributes

# Multivariate Normal Distribution

1-dimensional:  $x \in \mathbb{R}$

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2\sigma^2} (x - \mu)^2 \right]$$



$D$ -dimensional:  $\mathbf{x} \in \mathbb{R}^D$

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

# Multivariate Normal Distribution

□ Mahalanobis distance:  $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

measures the distance from  $\mathbf{x}$  to  $\boldsymbol{\mu}$  in terms of  $\boldsymbol{\Sigma}$   
normalizes for difference in variances and correlations

□ Bivariate:  $D = 2$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

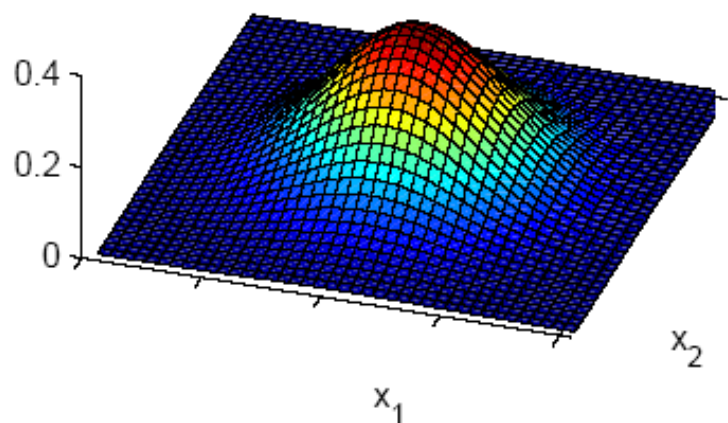
$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[ -\frac{1}{2(1-\rho^2)} (z_1^2 - 2\rho z_1 z_2 + z_2^2) \right]$$

$$z_i = \frac{x_i - \mu_i}{\sigma_i}$$

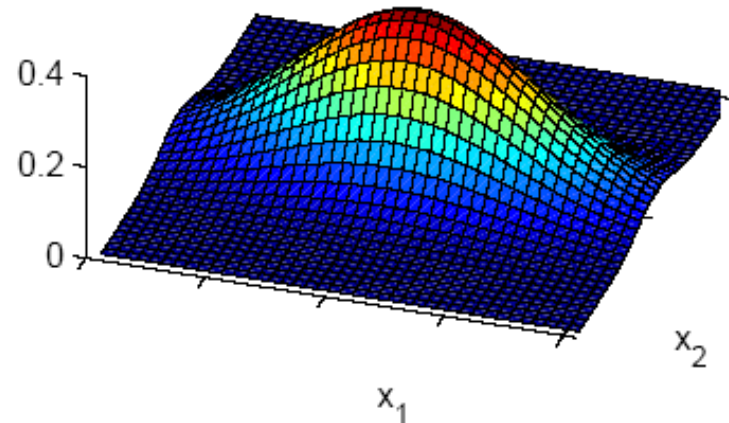


$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$$

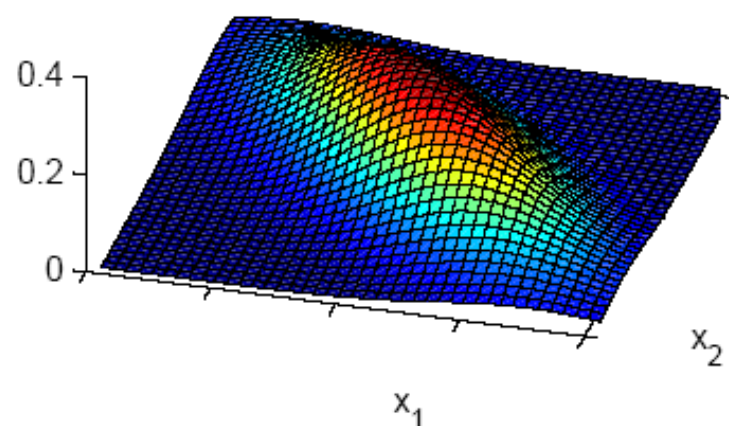
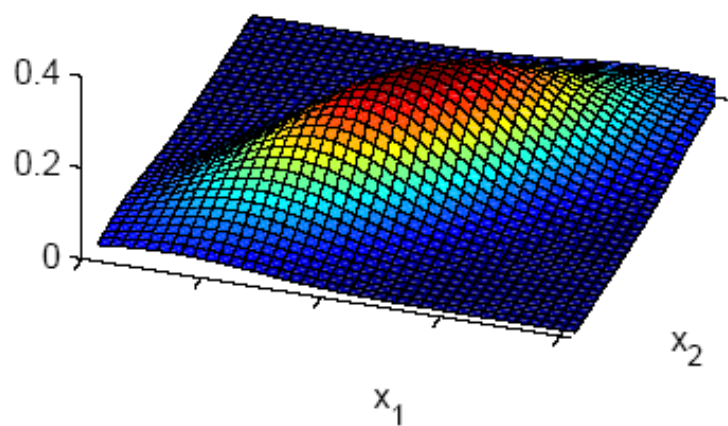
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$$



$$\text{Cov}(x_1, x_2) > 0$$

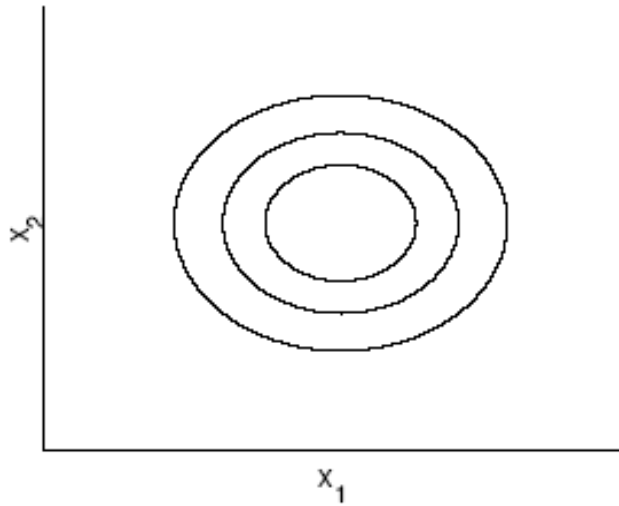


$$\text{Cov}(x_1, x_2) < 0$$

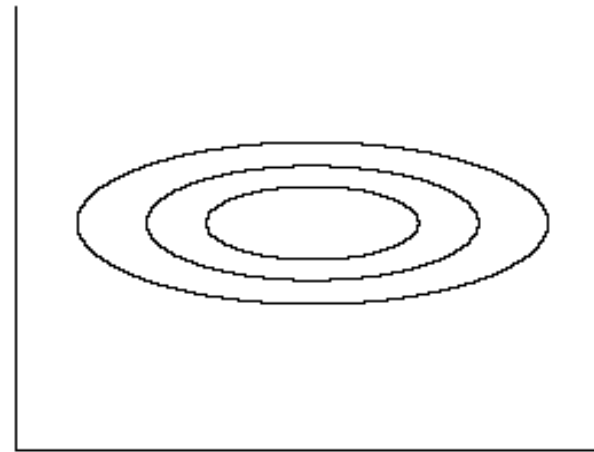


# Bivariate Normal

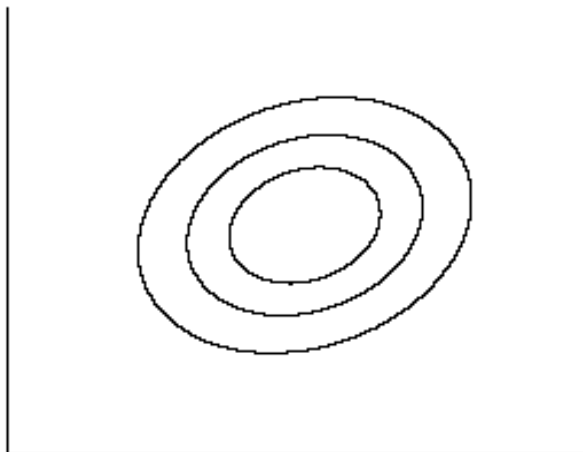
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$$



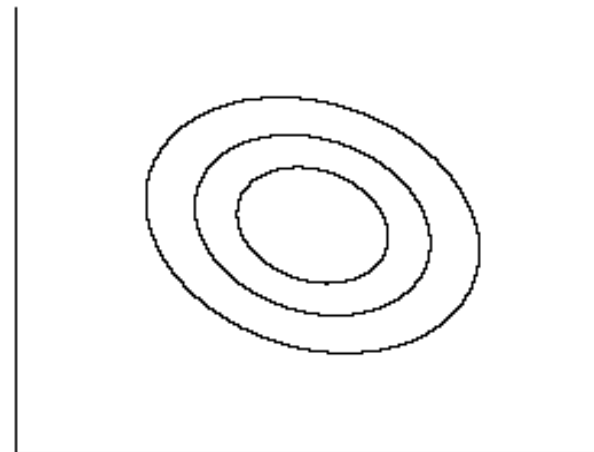
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$$



$$\text{Cov}(x_1, x_2) > 0$$



$$\text{Cov}(x_1, x_2) < 0$$



# Independent Inputs: Naive Bayes

- If  $x_i$  are independent
  - $x_i$  are uncorrelated
  - offdiagonals of  $\Sigma$  are 0
  - Mahalanobis distance reduces to weighted (by  $1/\sigma_d$ ) Euclidean distance:

$$p(\mathbf{x}) = \prod_{d=1}^D p_d(x_d) = \frac{1}{(2\pi)^{D/2} \prod_{d=1}^D \sigma_d} \exp \left[ -\frac{1}{2} \sum_{d=1}^D \left( \frac{x_d - \mu_d}{\sigma_d} \right)^2 \right]$$

- If variances are also equal  $\Rightarrow$  Euclidean distance

# Revisit: Parametric Classification Ch 4.5

- Suppose classes  $C_i, \quad i = 1, \dots, K$
- Discriminative function for each class  $g_i(\mathbf{x})$
- For sample  $\mathbf{x}$  choose class  $k$ :

$$g_k(\mathbf{x}) = \max_i g_i(\mathbf{x})$$

- How to define discriminate function  $g_i(\mathbf{x})$  ?
  - „Dumb“: class probability  $P(C_i)$
  - Maximum Likelihood (ML)  $p(\mathbf{x}|C_i)$
  - Maximum a-posteriori (MAP)  $p(\mathbf{x}|C_i)P(C_i)$
  - Bayes  $-R(\alpha_i|\mathbf{x})$


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

# Parametric Classification

□ If  $(\mathbf{x}|C_i) \sim \mathcal{N}(\mathbf{x}_i, \mathbf{\Sigma}_i)$

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

□ Discriminant functions (see Ch. 3.4)

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x}|C_i) + \log P(C_i) \\ &= -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{\Sigma}_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \\ &\quad + \log P(C_i) \end{aligned}$$

# Estimation of Parameters

$$\mathbb{1} \{ \mathbf{x}_k \in C_i \} = \begin{cases} 1 & , \mathbf{x}_k \in C_i \\ 0 & , \text{else} \end{cases}, \quad N_i := \sum_{k=1}^N \mathbb{1} \{ \mathbf{x}_k \in C_i \}$$

$$\hat{P}(C_i) = \frac{N_i}{N}$$

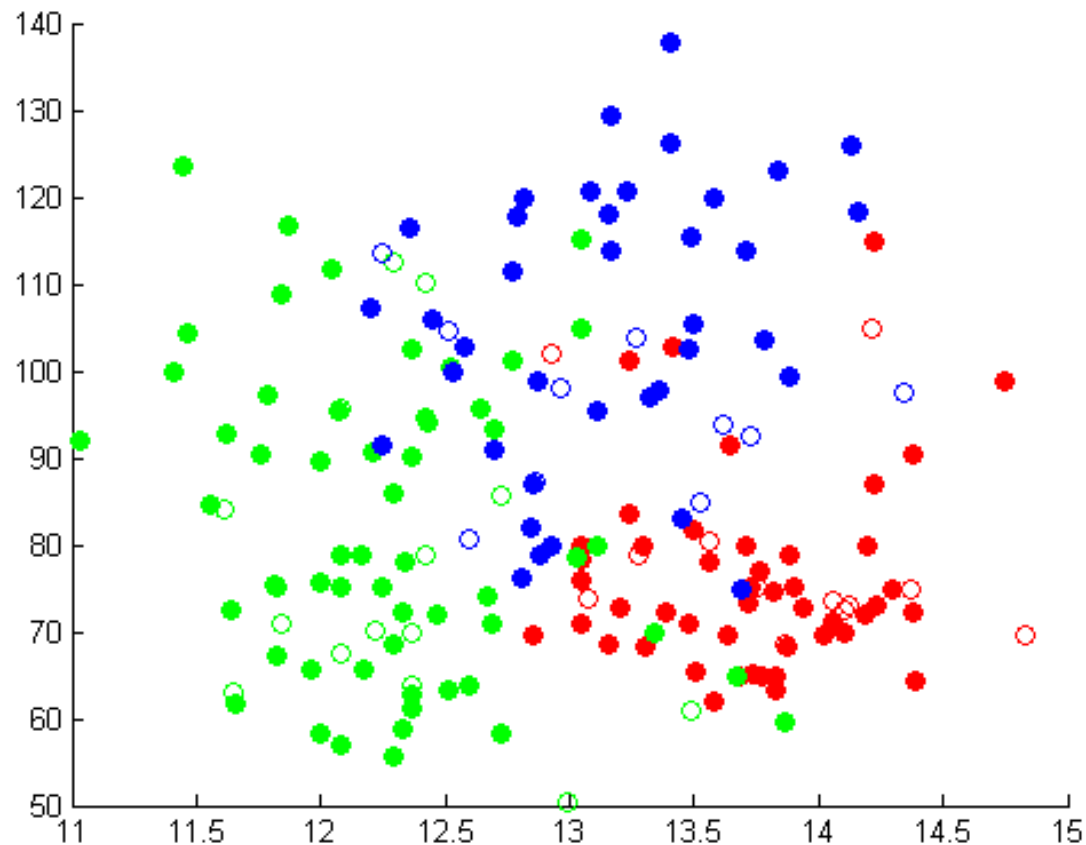
$$\mathbf{m}_i = \frac{1}{N_i} \sum_{k=1}^N \mathbb{1} \{ \mathbf{x}_k \in C_i \} \mathbf{x}_k$$

$$\mathbf{S}_i = \frac{1}{N_i} \sum_{k=1}^N \mathbb{1} \{ \mathbf{x}_k \in C_i \} (\mathbf{x}_k - \mathbf{m}_i)(\mathbf{x}_k - \mathbf{m}_i)^T$$

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

# Parametric Classification

- 3 classes
- Training and Test Data



# Case 1: Different $\mathbf{S}_i$

Quadratic discriminant

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

$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

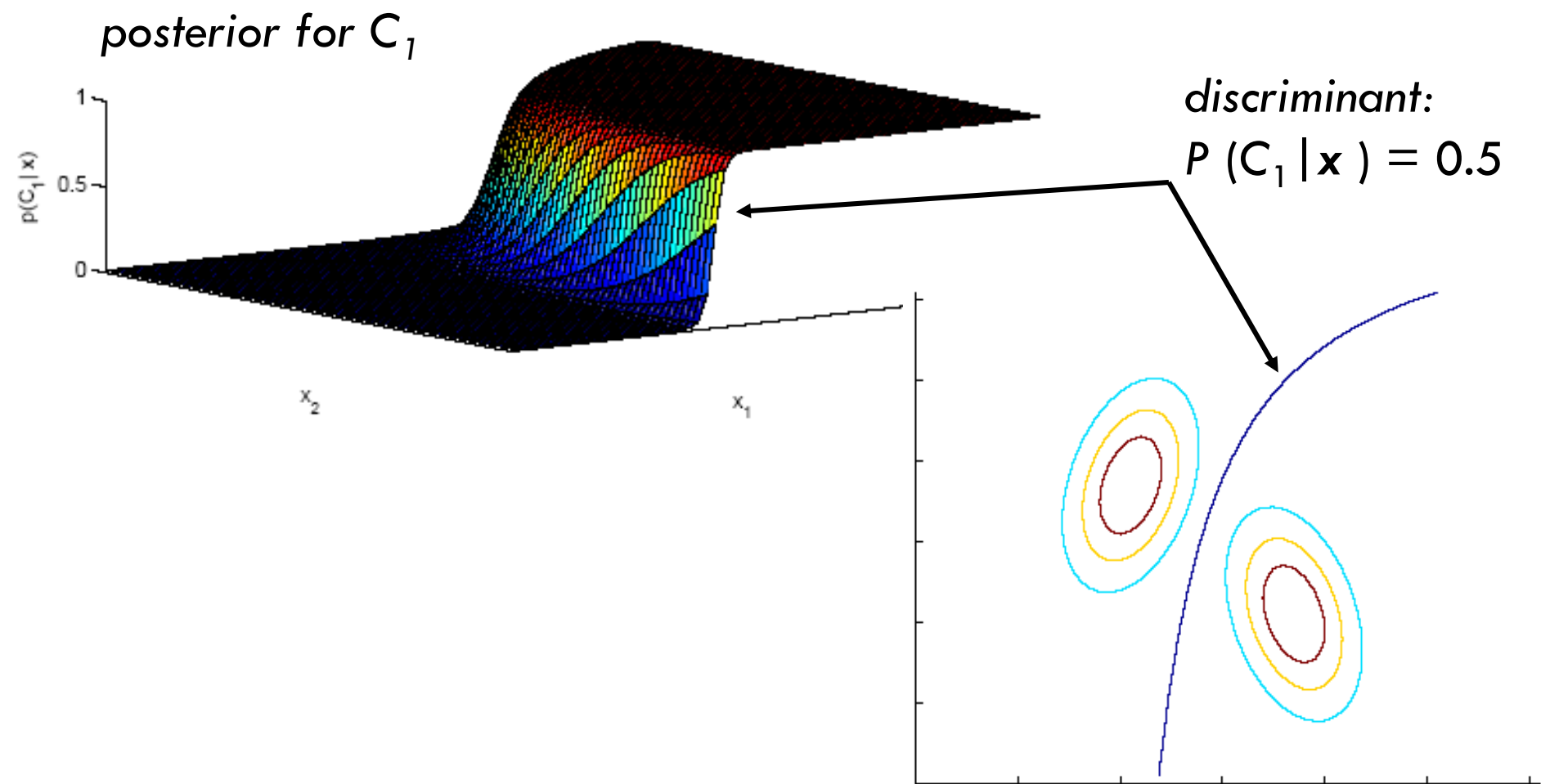
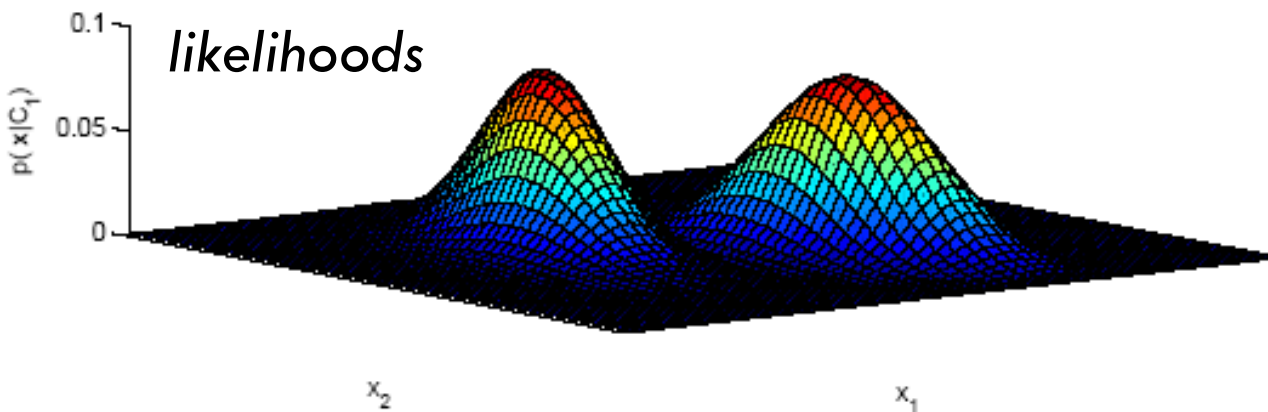
where

$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$$

$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$$

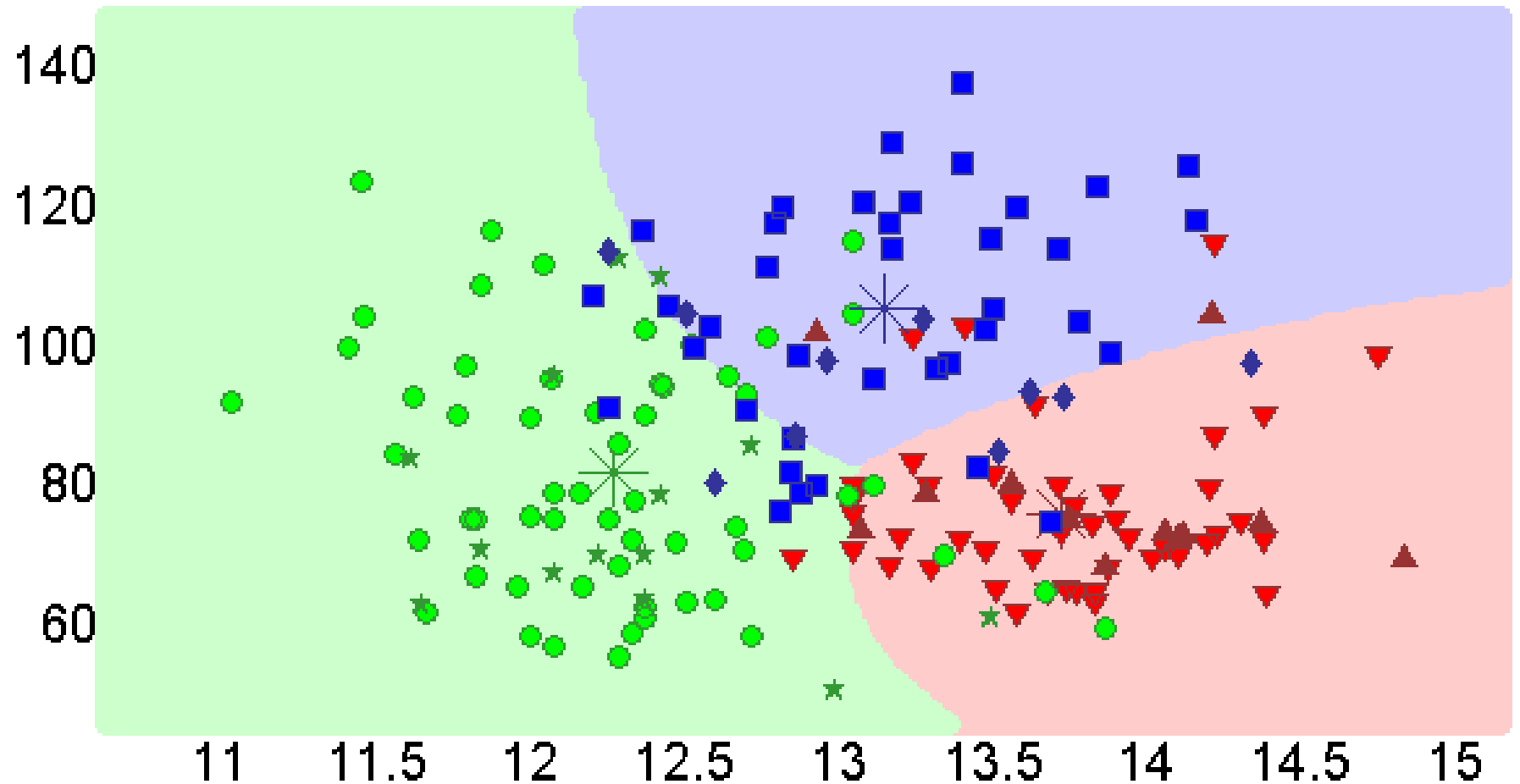
$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}_i| + \log \hat{P}(C_i)$$





# Case 1: Different $S_i$

type 1, 83.10/72.22 percent correct (train/test)



# Case 2: Common Covariance Matrix $\mathbf{S}$

- Shared common sample covariance  $\mathbf{S}$

$$\mathbf{S} = \sum_i \hat{P}(C_i) \mathbf{S}_i$$

- Discriminant 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

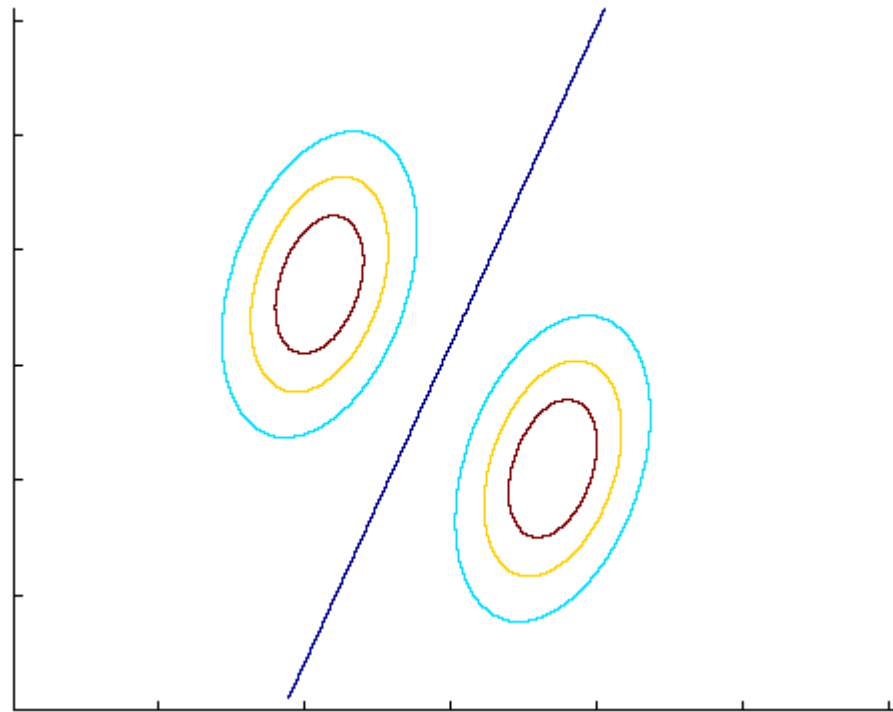
$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where

$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i$$

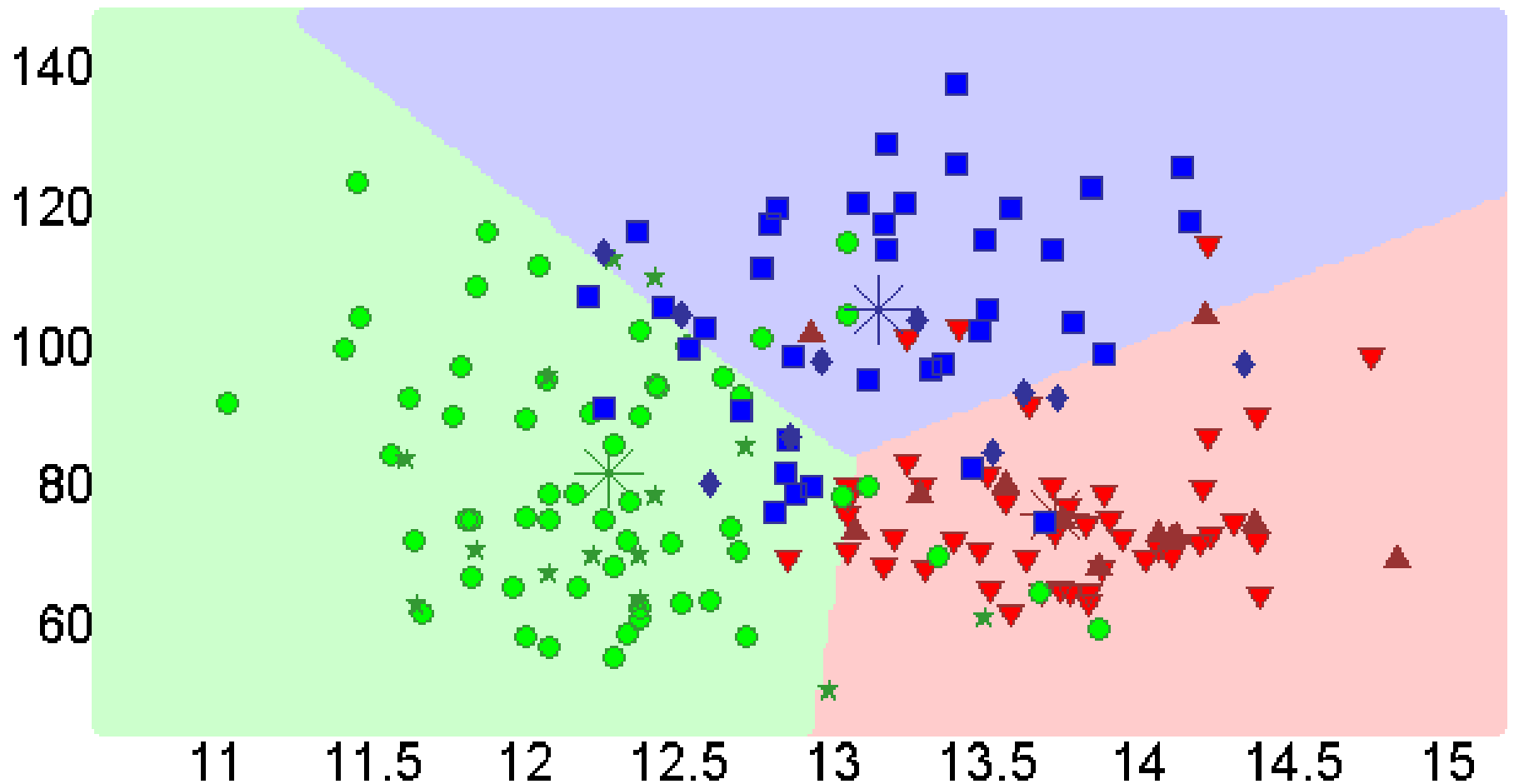
$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}(C_i)$$

# Case 2: Common Covariance Matrix $\mathbf{S}$



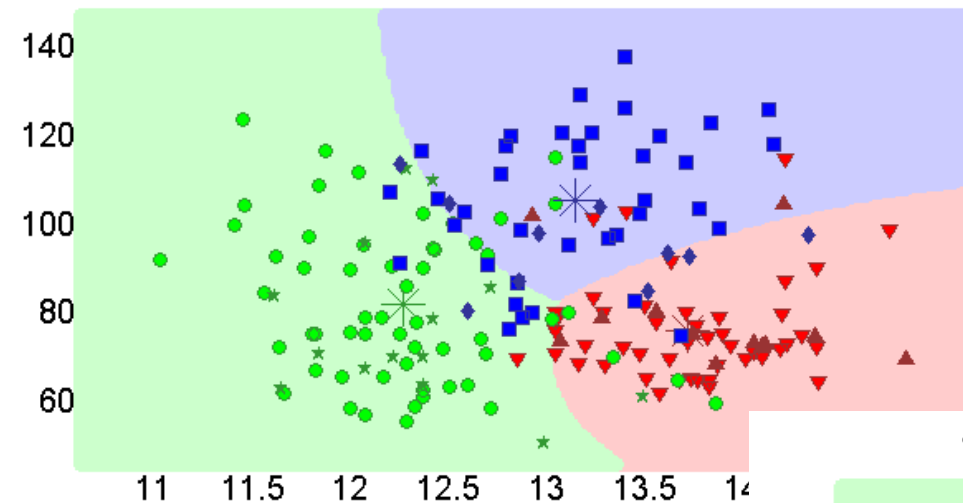
# Case 2: Common Covariance Matrix $S$

type 2, 82.39/66.67 percent correct (train/test)

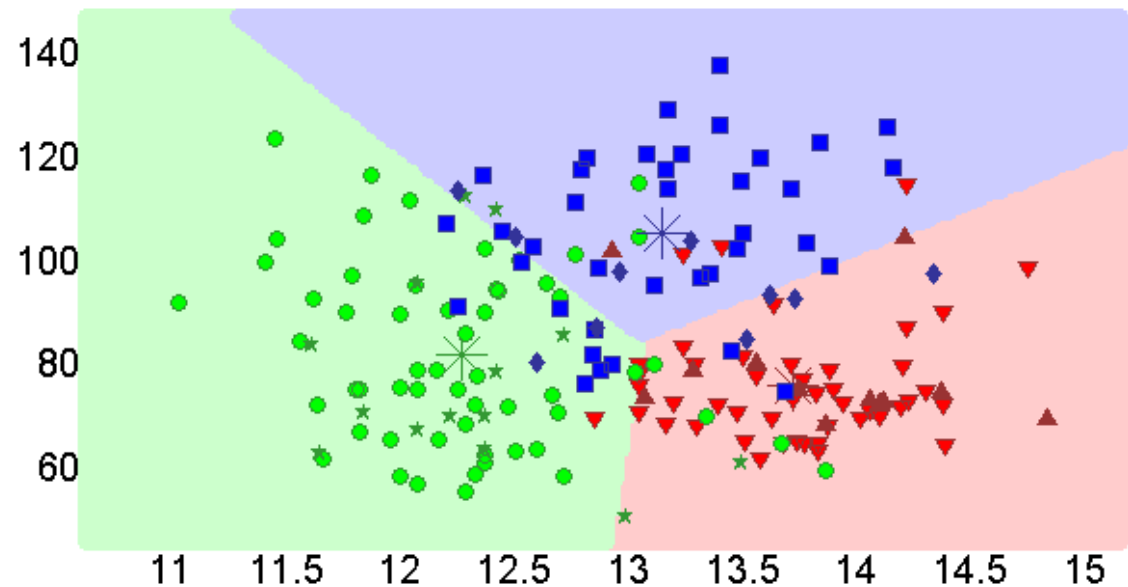


# Case 2: Common Covariance Matrix $\Sigma$

type 1, 83.10/72.22 percent correct (train/test)



type 2, 82.39/66.67 percent correct (train/test)



# Case 3: Diagonal S

If  $x_d$   $d = 1, \dots, D$  are independent:

- $p(\mathbf{x}|C_i) = \prod_{d=1}^D p(x_d|C_i)$  (Naive Bayes' assumption)

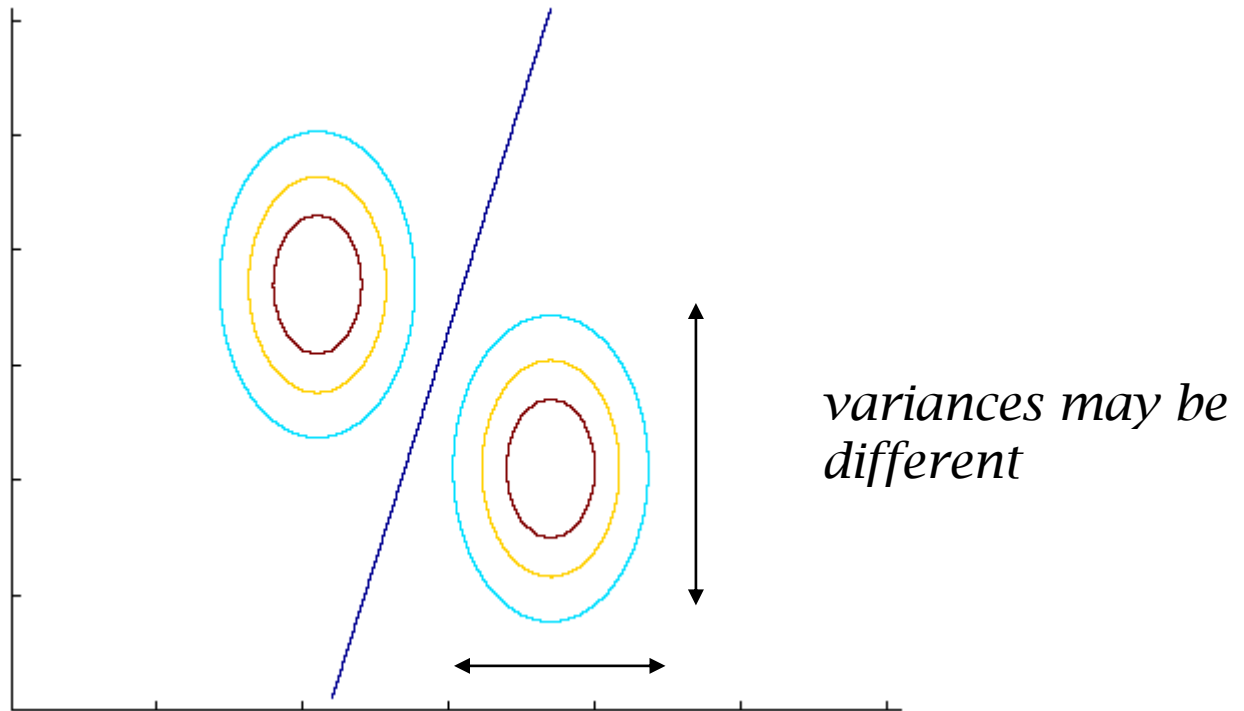
- Covariance matrix is diagonal

$$g_i(\mathbf{x}) = -\frac{1}{2} \sum_{d=1}^D \left( \frac{x_d - m_{id}}{s_d} \right)^2 + \log \hat{P}(C_i)$$

Classify based on

weighted Euclidean distance (in  $s_d$  units) to the nearest mean

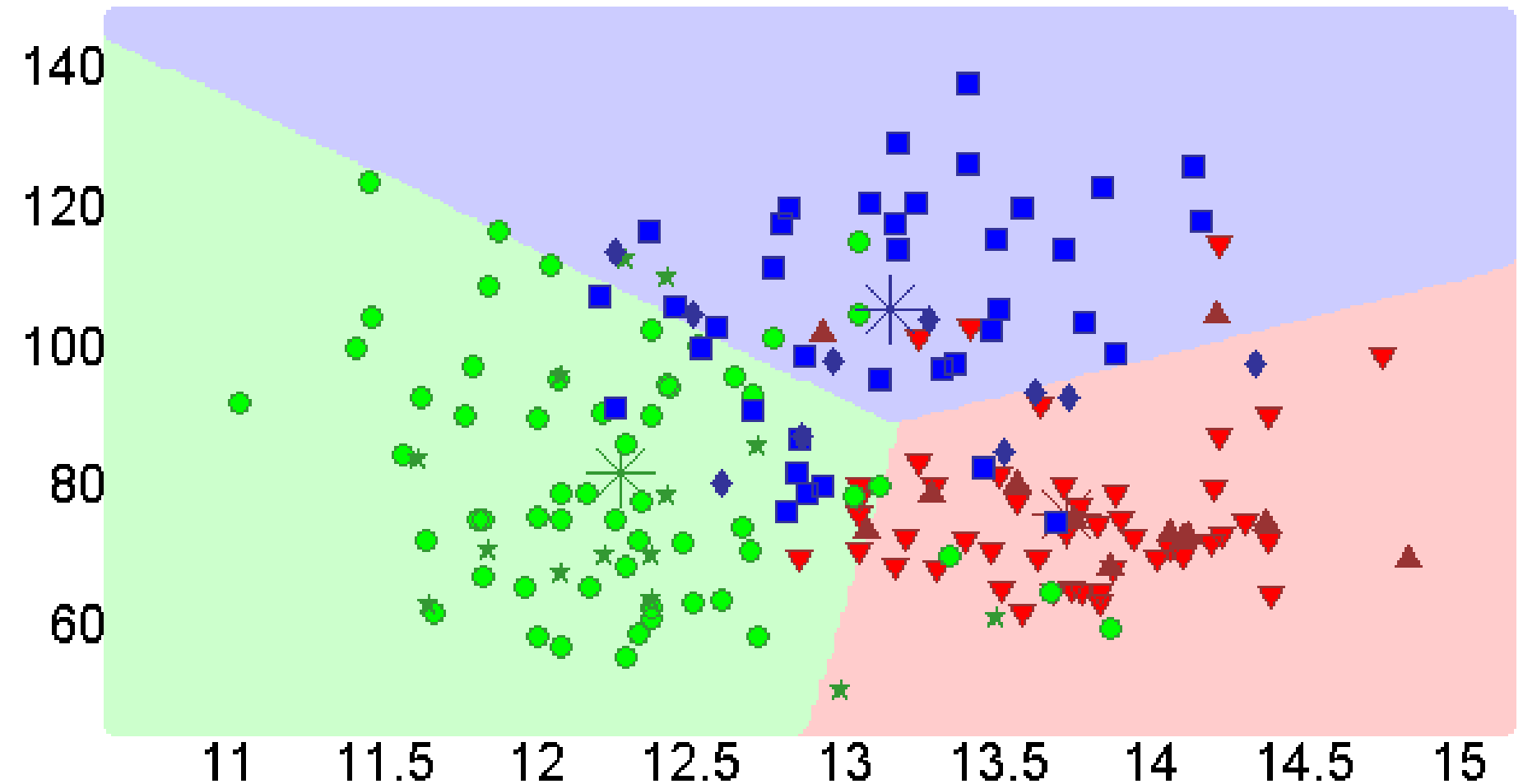
# Case 3: Diagonal S





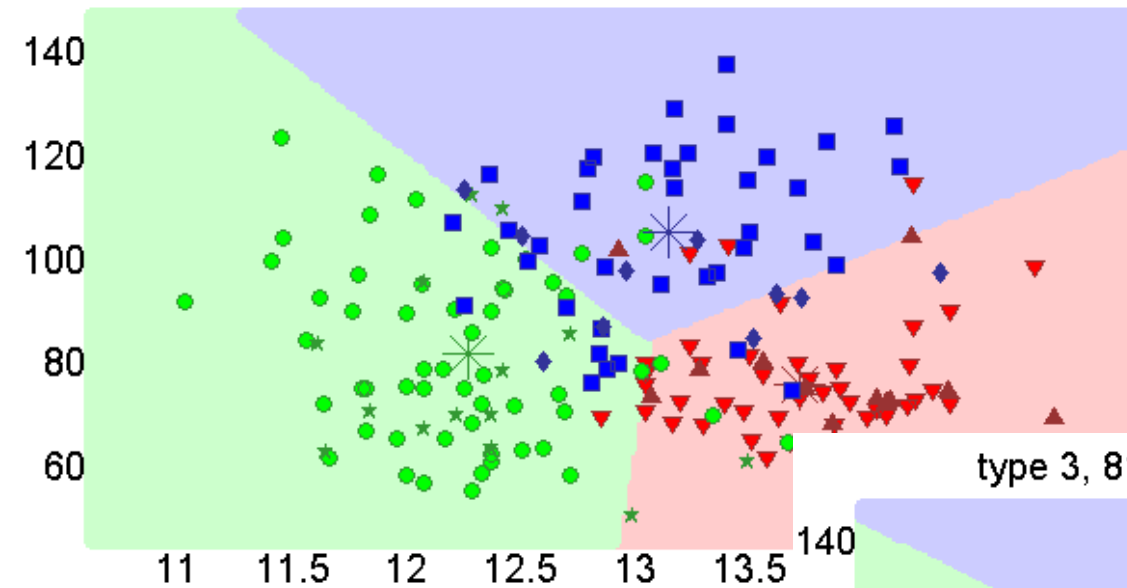
# Case 3: Diagonal S

type 3, 81.69/66.67 percent correct (train/test)

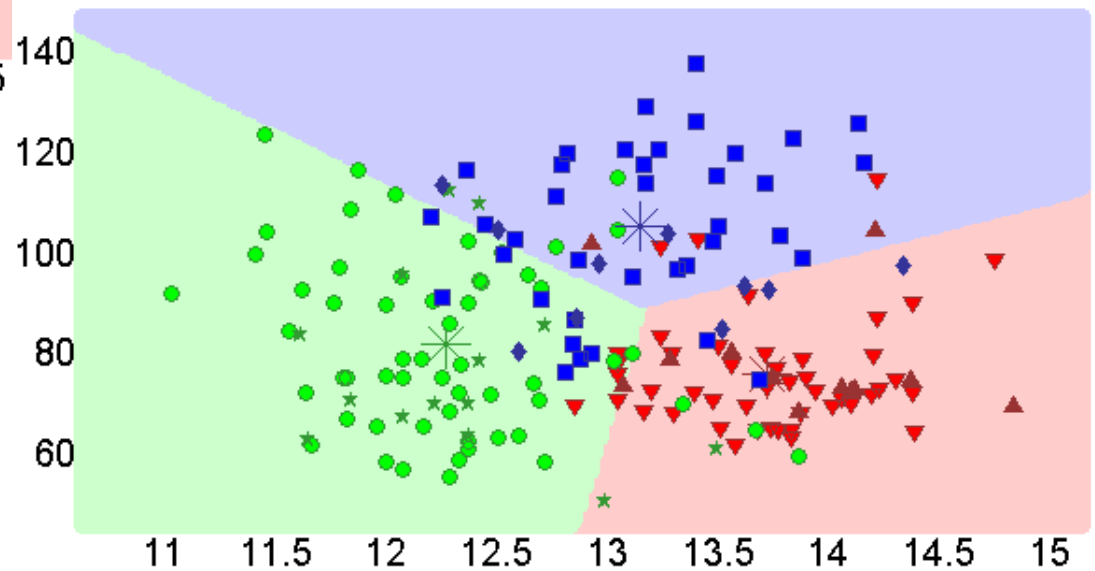


# Case 3: Diagonal S

type 2, 82.39/66.67 percent correct (train/test)



type 3, 81.69/66.67 percent correct (train/test)



# Case 4: Diagonal $\mathbf{S}$ , equal variances

## **Nearest mean classifier:**

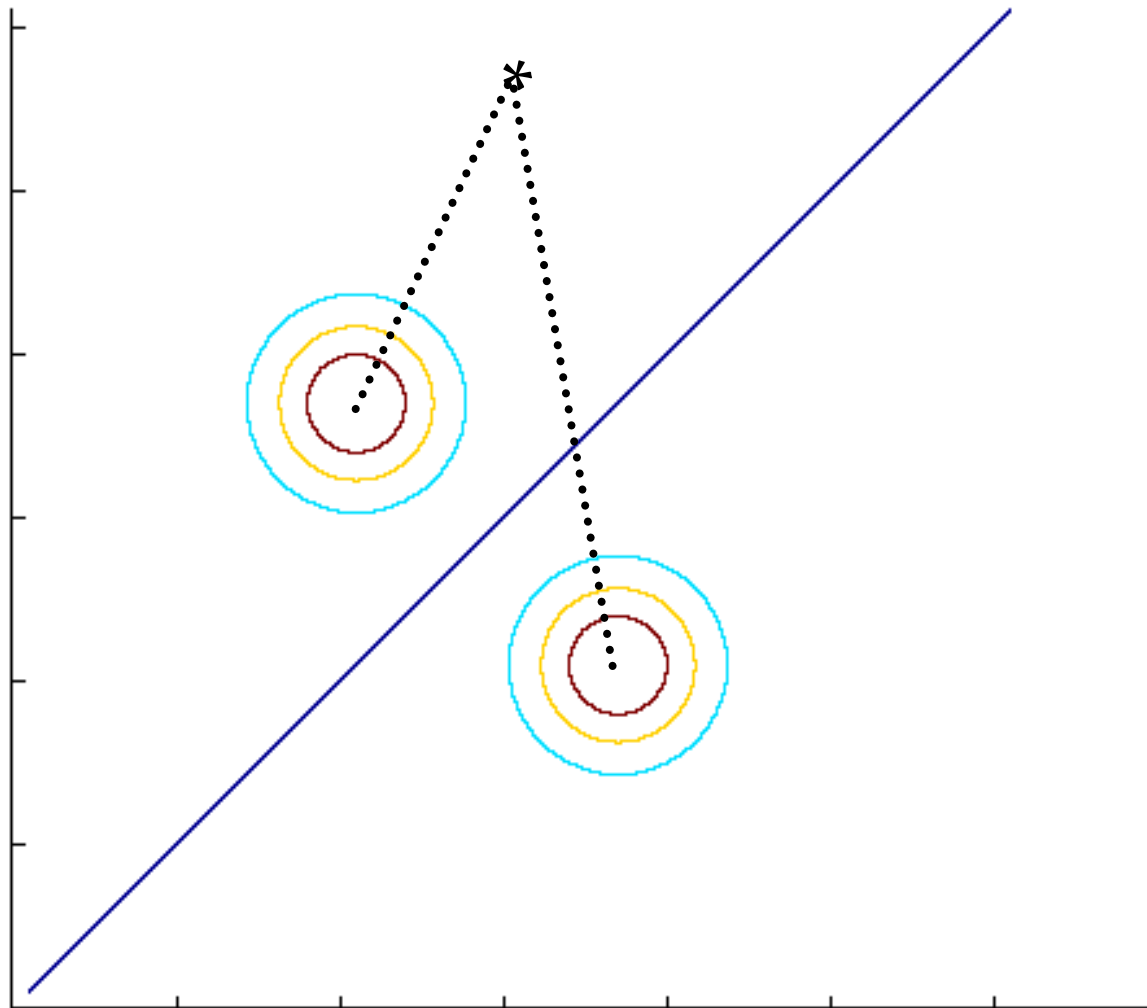
Classify based on Euclidean distance to the nearest mean

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

Each mean can be considered a prototype or template,  
hence this is **template matching**

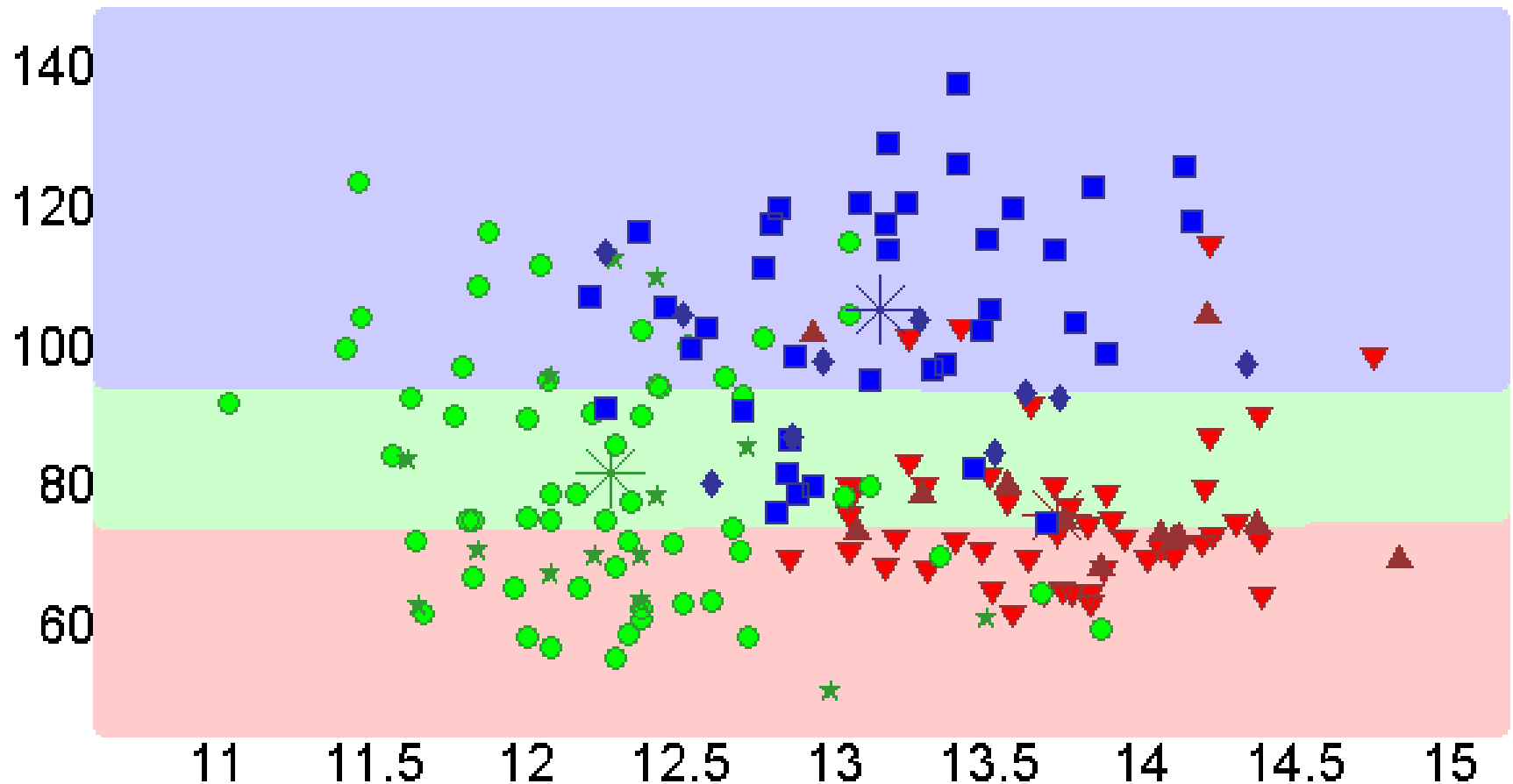
# Case 4: Diagonal $\mathbf{S}$ , equal variances

*Assign the closest mean value*



# Case 4: Diagonal $\mathbf{S}$ , equal variances

type 4, 52.82/44.44 percent correct (train/test)



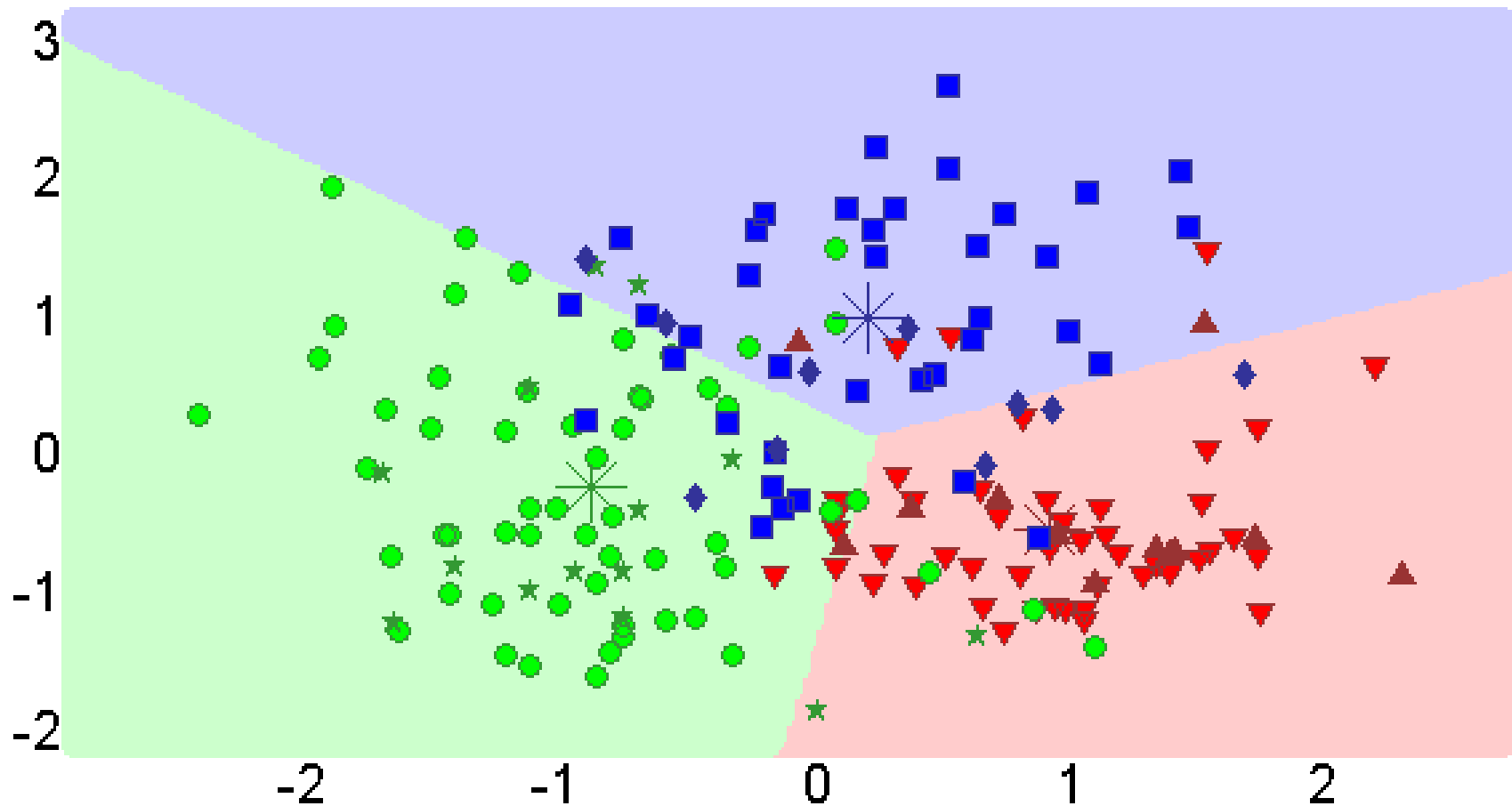
Why does it look so bad?

> No preprocessing: different variances (see axis)

# Case 4: Diagonal $\mathbf{S}$ , equal variances

Result **with** preprocessing by z-normalization

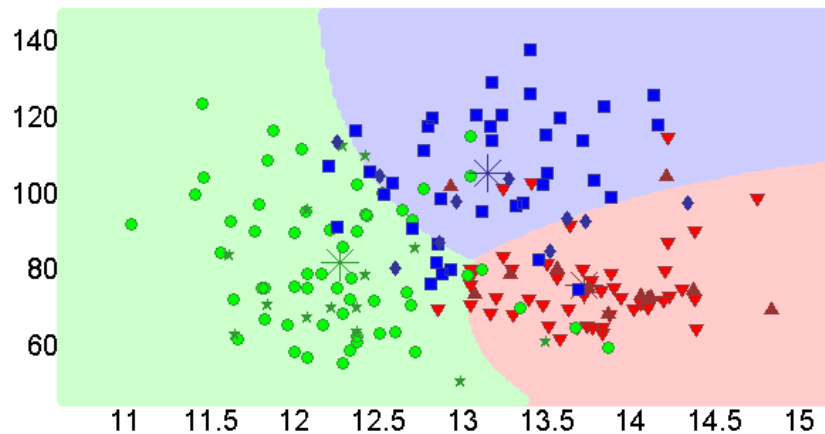
type 4, 81.69/66.67 percent correct (train/test)



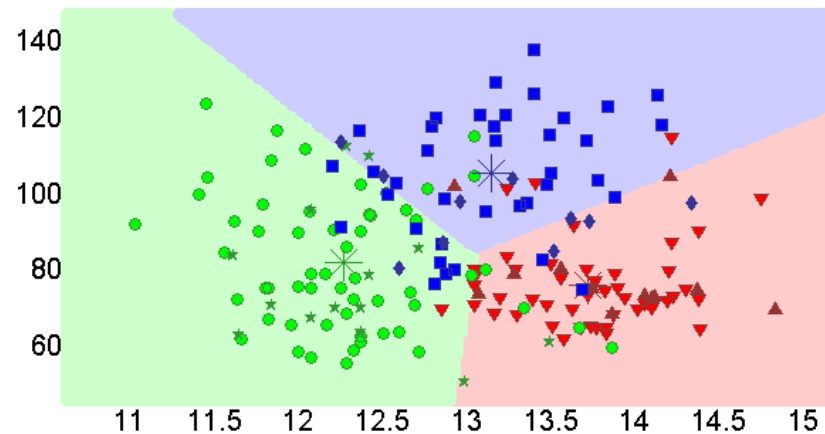
# Parametric Classification

## Result **without** preprocessing by z-normalization

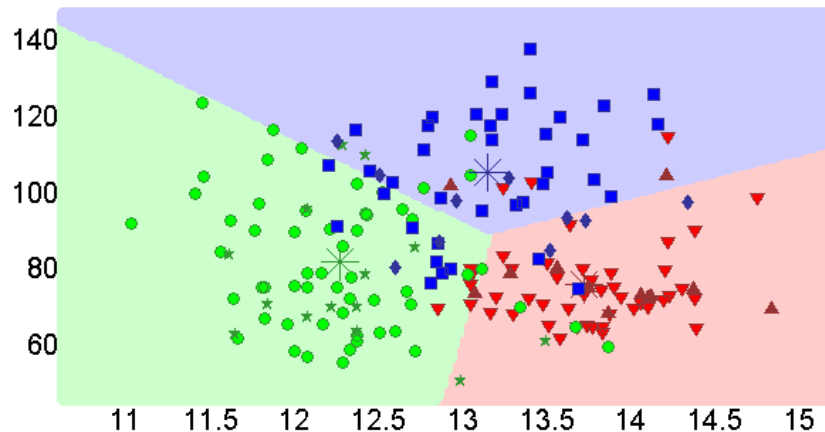
type 1, 83.10/72.22 percent correct (train/test)



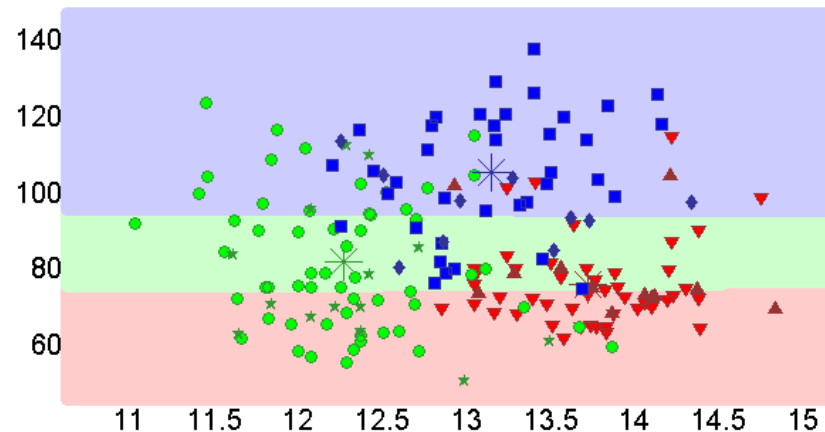
type 2, 82.39/66.67 percent correct (train/test)



type 3, 81.69/66.67 percent correct (train/test)



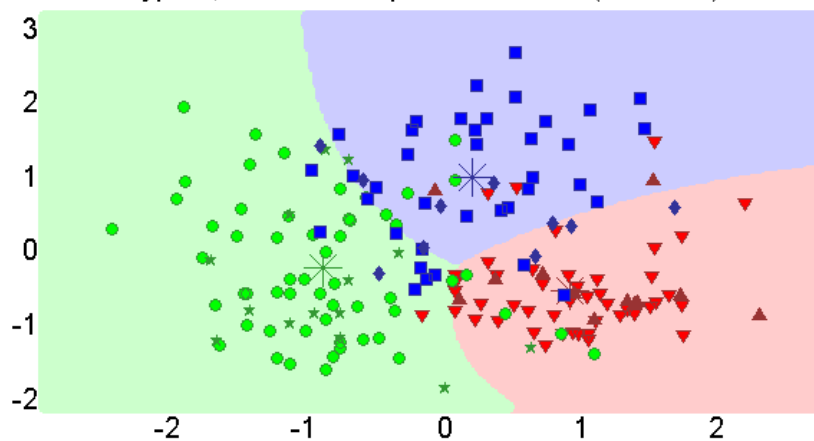
type 4, 52.82/44.44 percent correct (train/test)



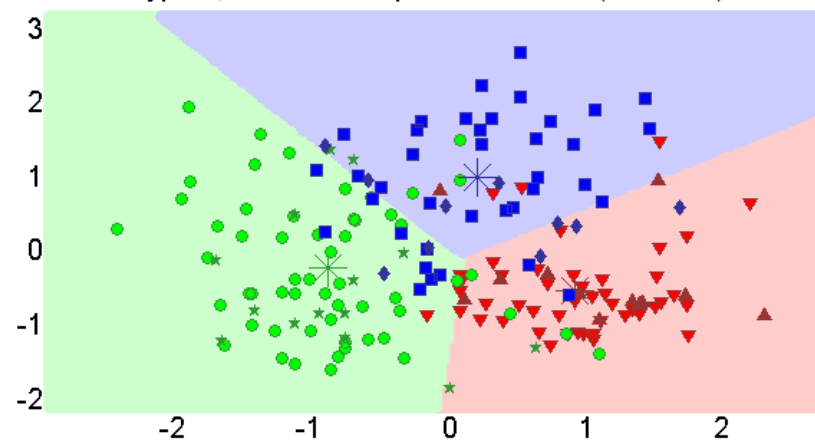
# Parametric Classification

## Result **with** preprocessing by z-normalization

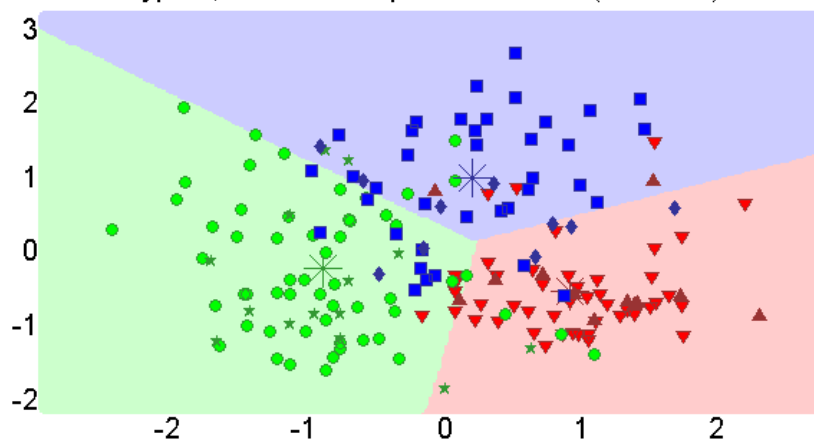
type 1, 83.10/72.22 percent correct (train/test)



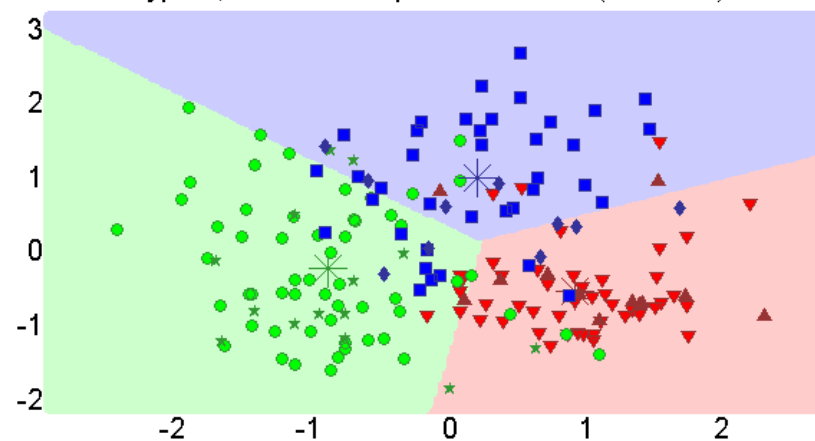
type 2, 82.39/66.67 percent correct (train/test)



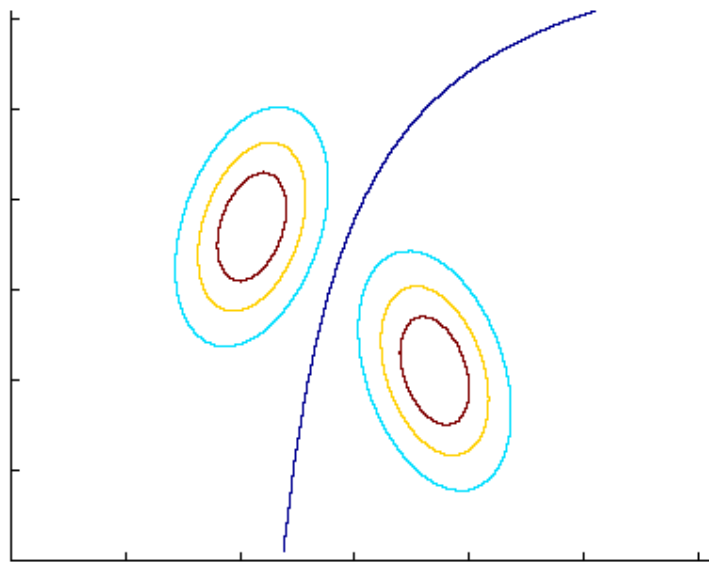
type 3, 81.69/66.67 percent correct (train/test)



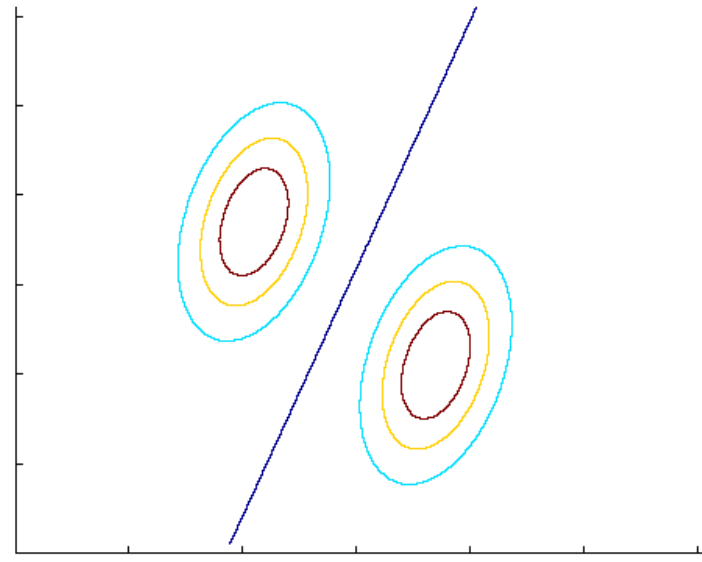
type 4, 81.69/66.67 percent correct (train/test)





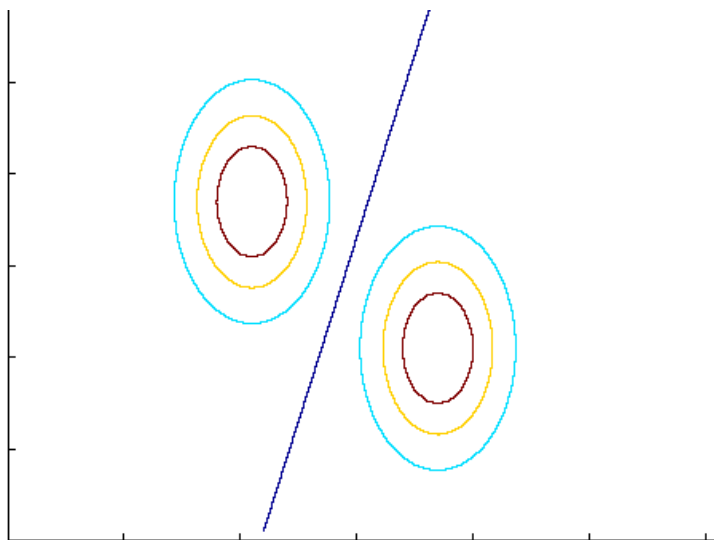


individual  $\mathbf{S}_i$

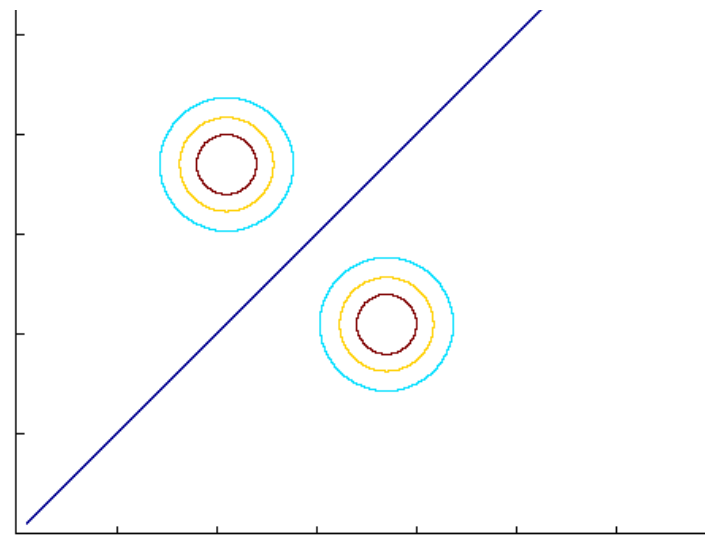


one  $\mathbf{S}$  for all

diagonal  $\mathbf{S}$



diagonal  $\mathbf{S}$  same variance



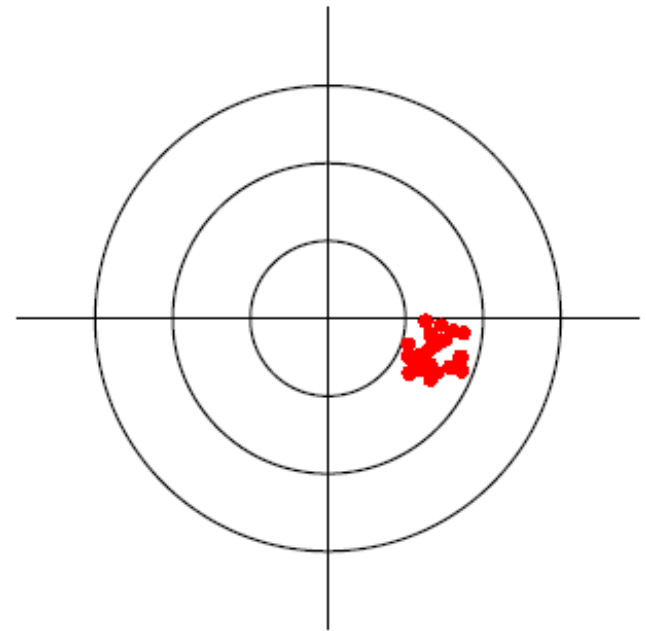
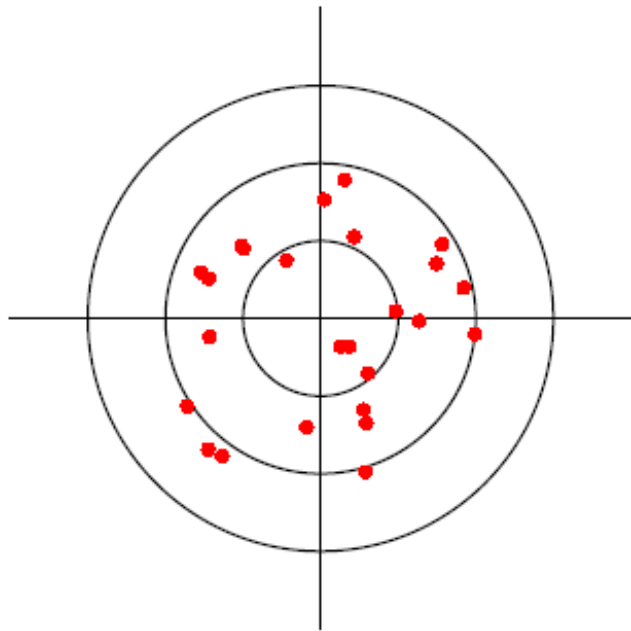
# Model Selection

<i>Assumption</i>	<i>Covariance matrix</i>	<i>No of parameters</i>
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis-aligned	$\mathbf{S}_i = \mathbf{S}$ , with $s_{ij} = 0$	$D$
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$D(D+1)/2$
Different, Hyperellipsoidal	$\mathbf{S}_i$	$K D(D+1)/2$

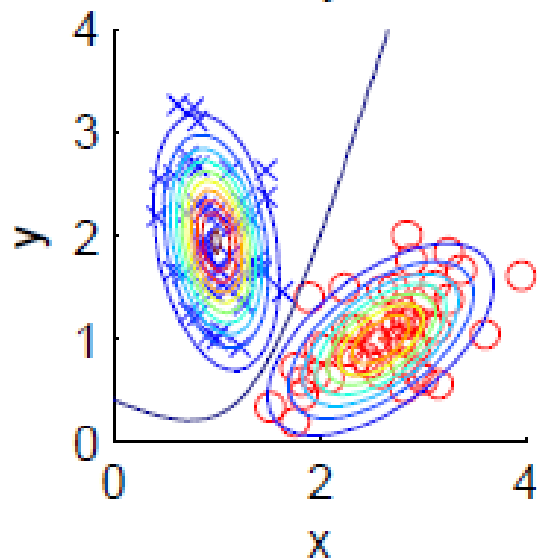
- As we increase complexity (less restricted  $\mathbf{S}$ ):  
bias decreases and variance increases
- Assume simple models (allow some bias) to control  
variance (regularization)

# Model Selection

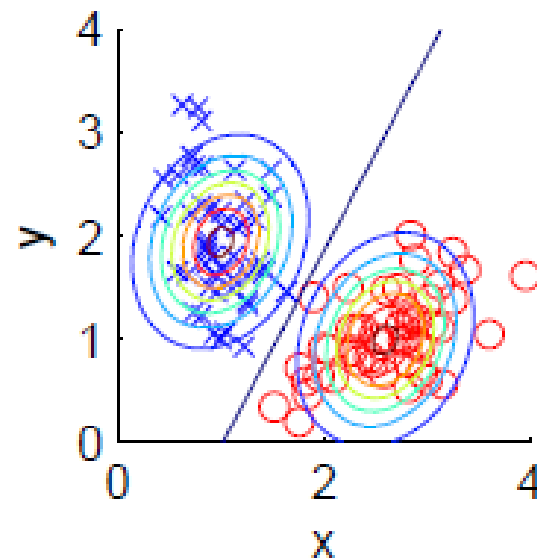
## Variance vs. Bias



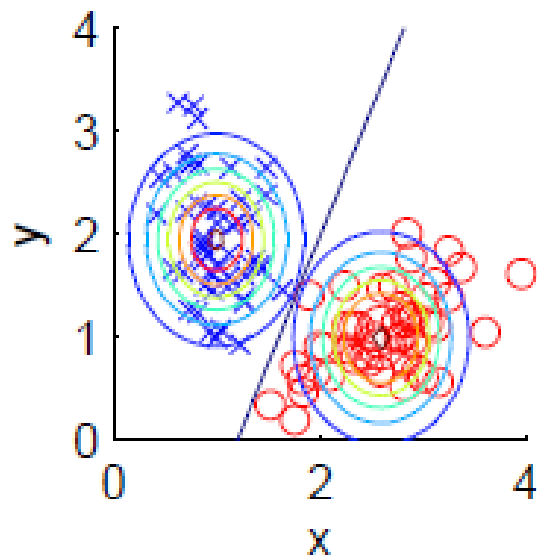
Arbitrary covar.



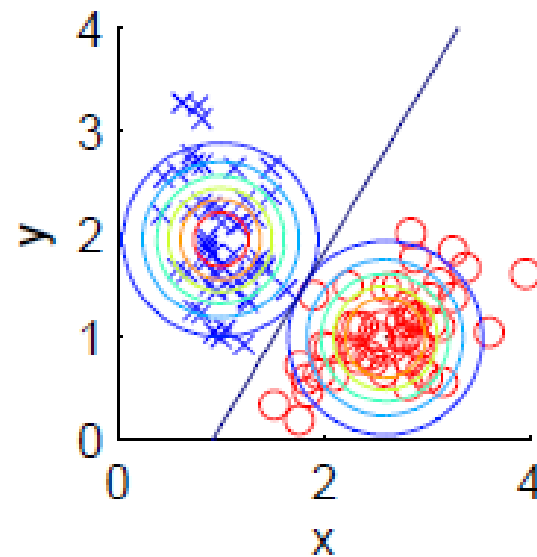
Shared covar.



Diag. covar.



Equal var.



# Model Selection

<i>Assumption</i>	<i>Covariance matrix</i>	<i>No of parameters</i>
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis-aligned	$\mathbf{S}_i = \mathbf{S}$ , with $s_{ij} = 0$	$D$
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$D(D+1)/2$
Different, Hyperellipsoidal	$\mathbf{S}_i$	$K D(D+1)/2$

- Weight the special cases:

$$\tilde{\mathbf{S}}_i = \alpha \sigma^2 \mathbf{I} + \beta \mathbf{S} + (1 - \alpha - \beta) \mathbf{S}_i, \quad \alpha, \beta \in [0, 1]$$

- Cross validation for selection

# General Considerations

prior knowledge of the data

- (in)dependence
- Preprocessing is important
- How to define model quality?  
Depends on application

# Overview of today

- 1) reminder and hint
- 2) multivariate data
  - ▣ descriptive statistics
  - ▣ classification
  - ▣ model selection
- 3) regression
  - ▣ 1D regression revisited
  - ▣ multivariate regression
  - ▣ model selection

# Revisit: Linear Model for Regression

- Consider data:  $(x_i, y_i)_{i=1}^N$
- linear model:  $y_i = w_0 + w_1 x_i + \epsilon_i, \quad i = 1, \dots, N$   
 $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$
- How to estimate parameters?

$$\min \frac{1}{2} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2$$
$$f(w_0, w_1) = \frac{1}{2} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2$$



# Revisit: Linear Model for Regression

$$y_i = w_0 + w_1 x_i + \epsilon_i, \quad i = 1, \dots, N, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\min \frac{1}{2} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2$$

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{X} := \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \\ 1 & x_N \end{pmatrix}, \quad \mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

How to rewrite error in matrix-vector representation?

# Revisit: Linear Model for Regression

$$y_i = w_0 + w_1 x_i + \epsilon_i, \quad i = 1, \dots, N, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\min \frac{1}{2} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2$$

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{X} := \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \\ 1 & x_N \end{pmatrix}, \quad \mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

$$\min \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \quad \text{What are the estimated parameters?}$$

# Revisit: Linear Model for Regression

$$y_i = w_0 + w_1 x_i + \epsilon_i, \quad i = 1, \dots, N, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\min \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

$$\begin{aligned} f(\mathbf{w}) &= \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2} (\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}) \end{aligned}$$

$$\nabla f(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} \stackrel{!}{=} \mathbf{0}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

# Revisit: Linear Model for Regression

$$y = Xw + \epsilon$$

$$\hat{w} = (X^T X)^{-1} X^T y$$

How is the estimator used to predict?

$$\hat{y} = X\hat{w}$$

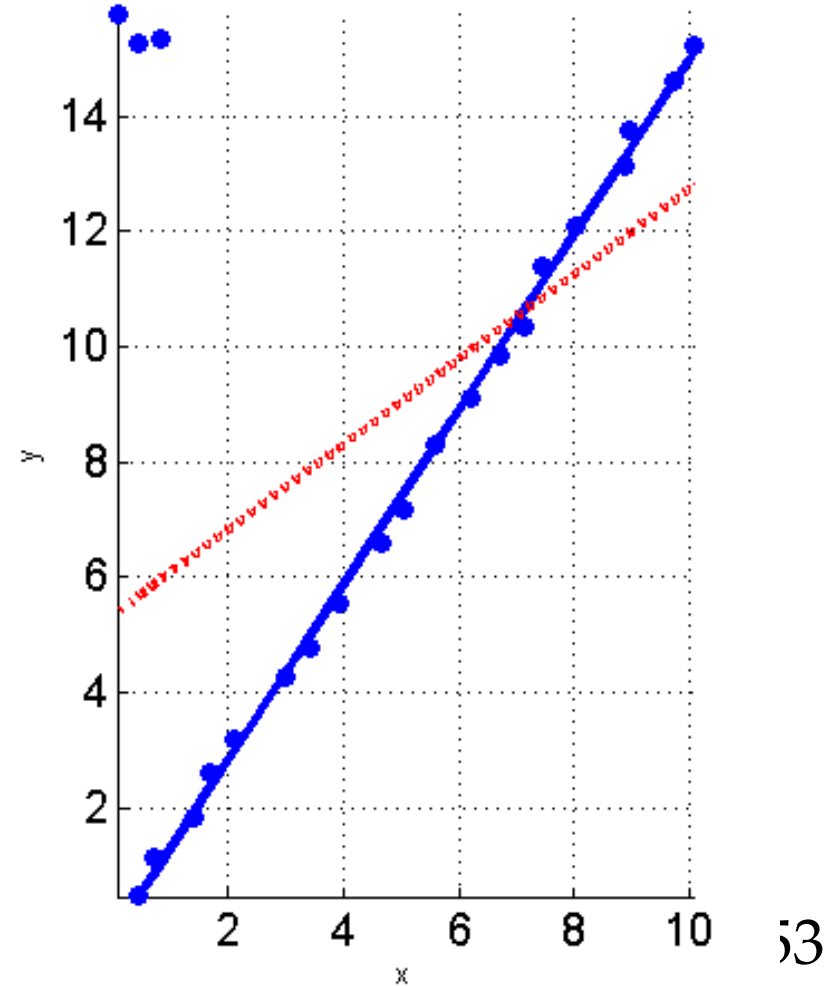
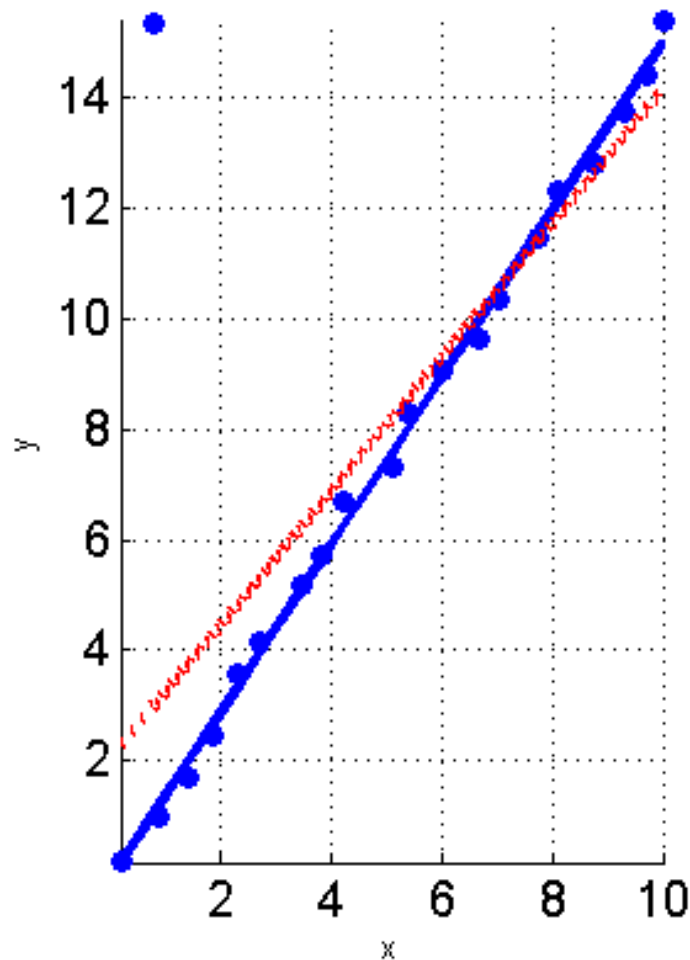
$$\hat{y} = \underbrace{X(X^T X)^{-1} X^T}_H y$$

- H is called the **hat matrix**
- $h_{ii}$  are **leverages**  $[\frac{1}{n}, 1]$
- Encodes influence of input points in y

# Revisit: Linear Model for Regression

## Outliers

Regression line without (blue) vs. **with** outliers (red)

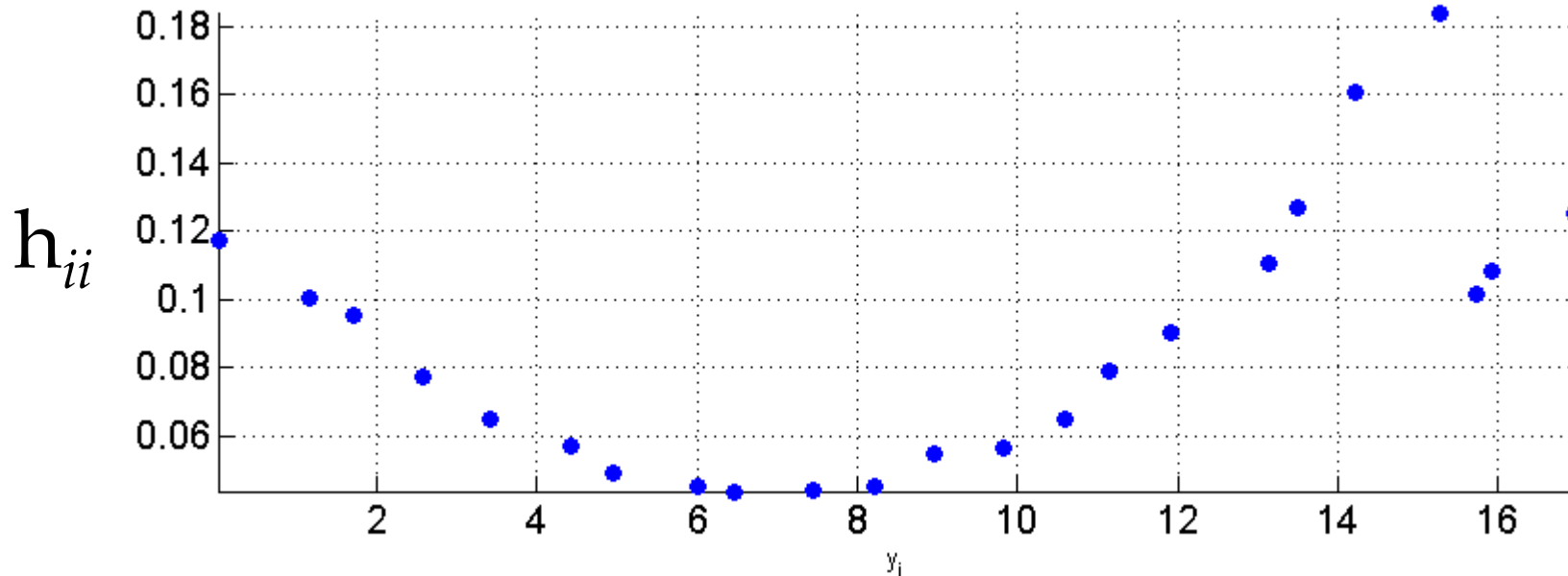
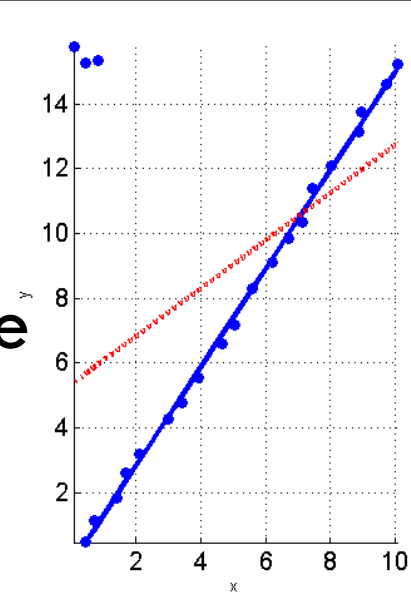


# Revisit: Linear Model for Regression

## Diagnostic Insights

Influence as **leverage** of points:

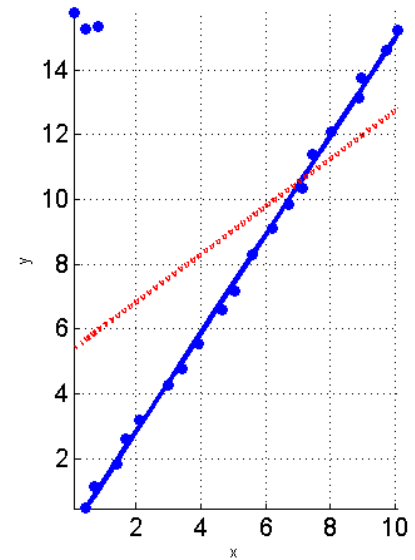
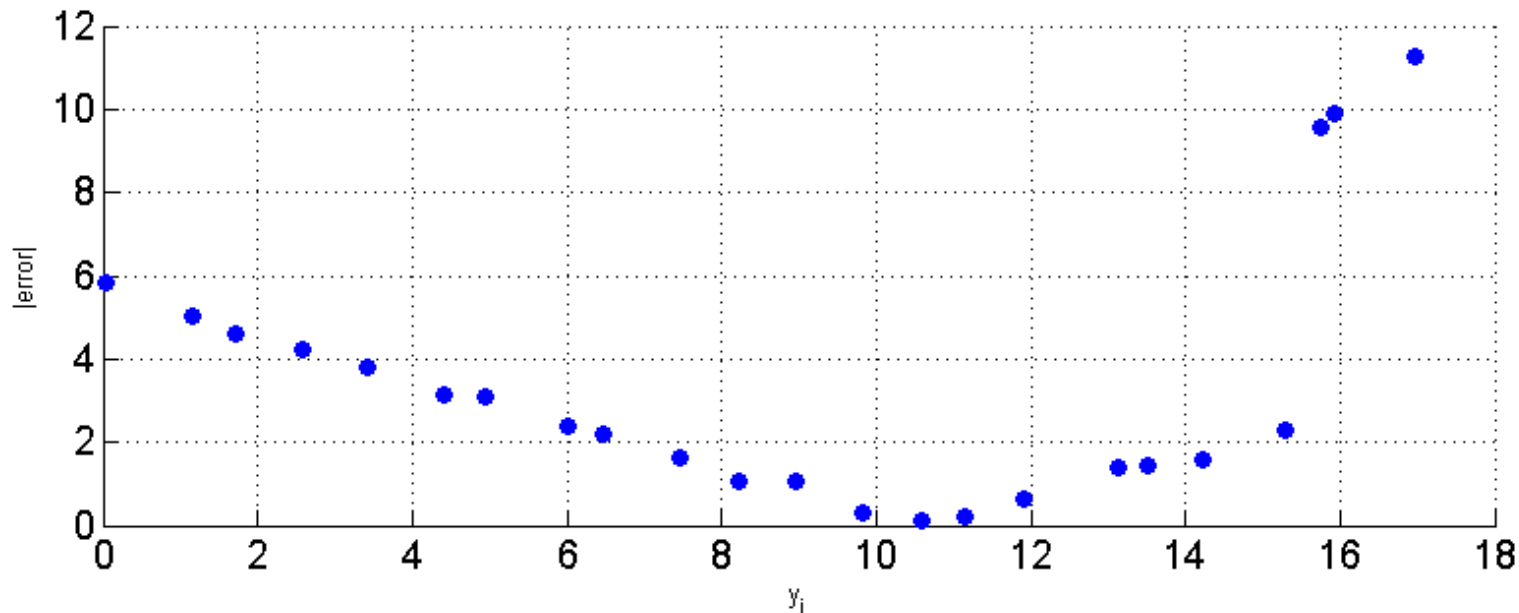
How much does point  $i$  influence the estimate



# Revisit: Linear Model for Regression

## Diagnostic Insights

Absolute error  $|\hat{y}_i - y_i|$

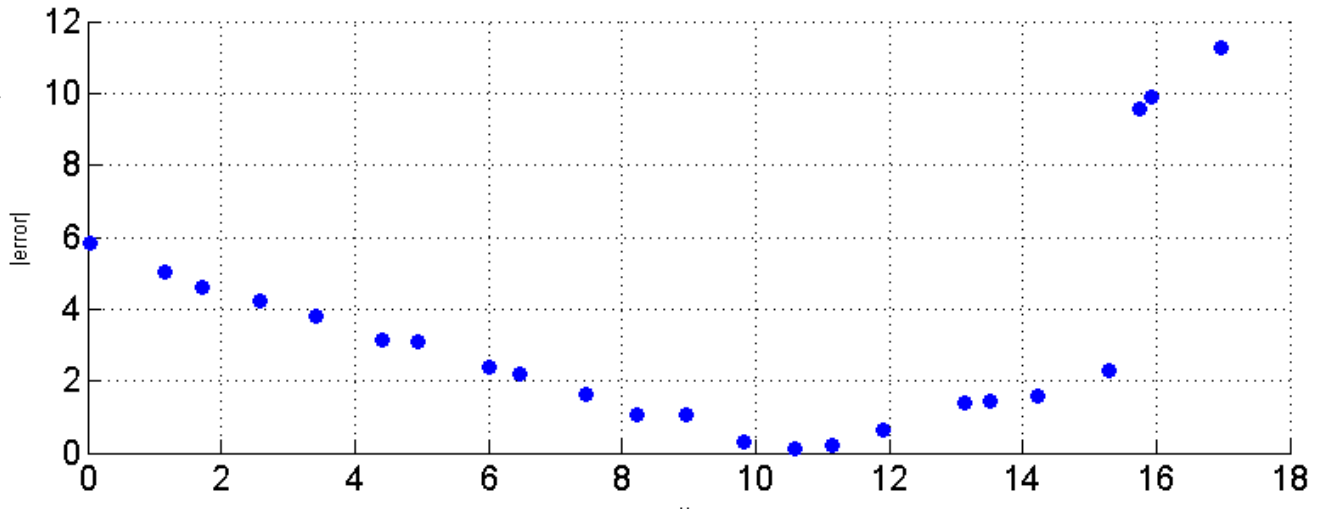


# Revisit: Linear Model for Regression

## Diagnostic Insights

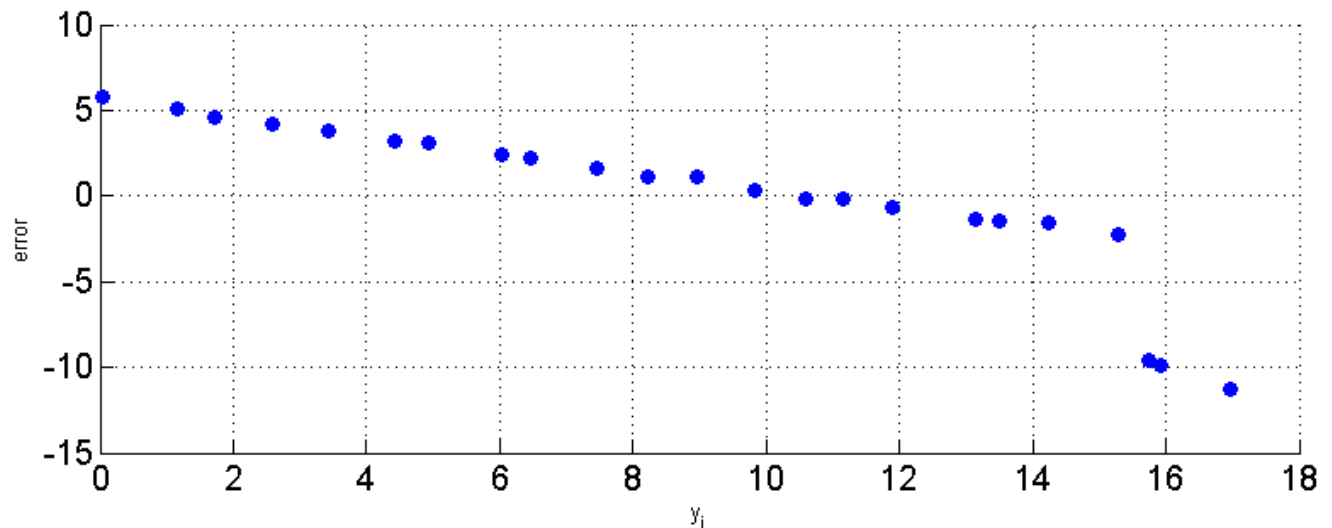
Absolute error

$$|\hat{y}_i - y_i|$$



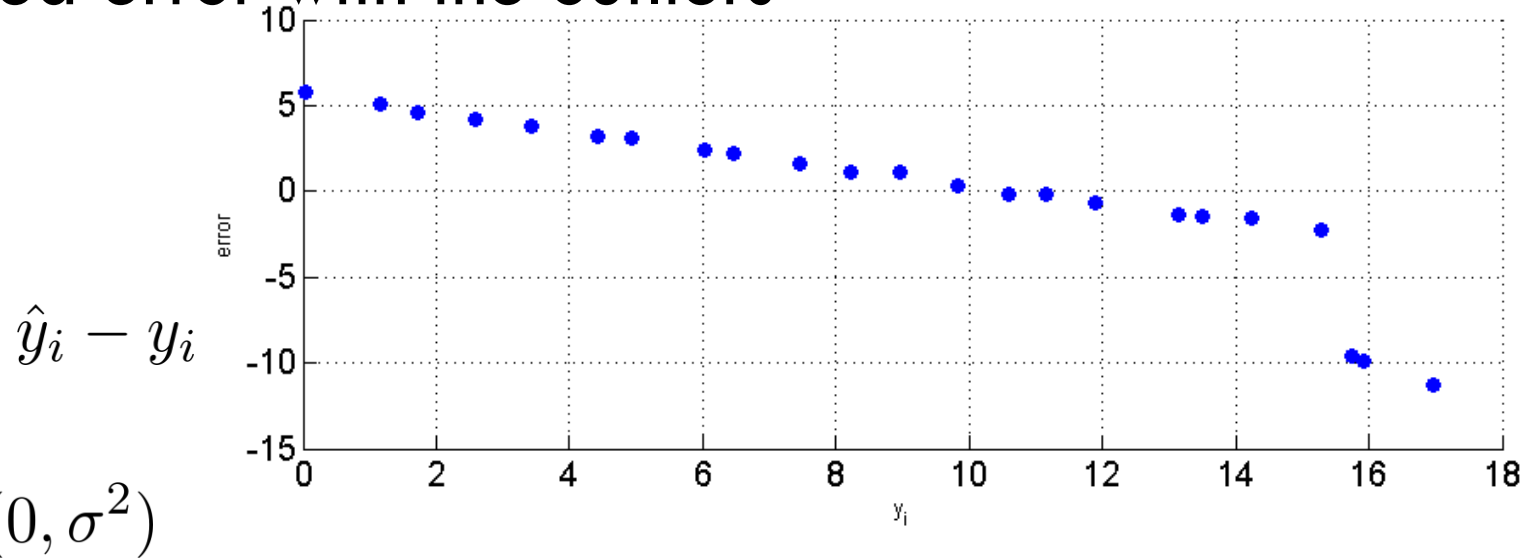
error with sign

$$\hat{y}_i - y_i$$

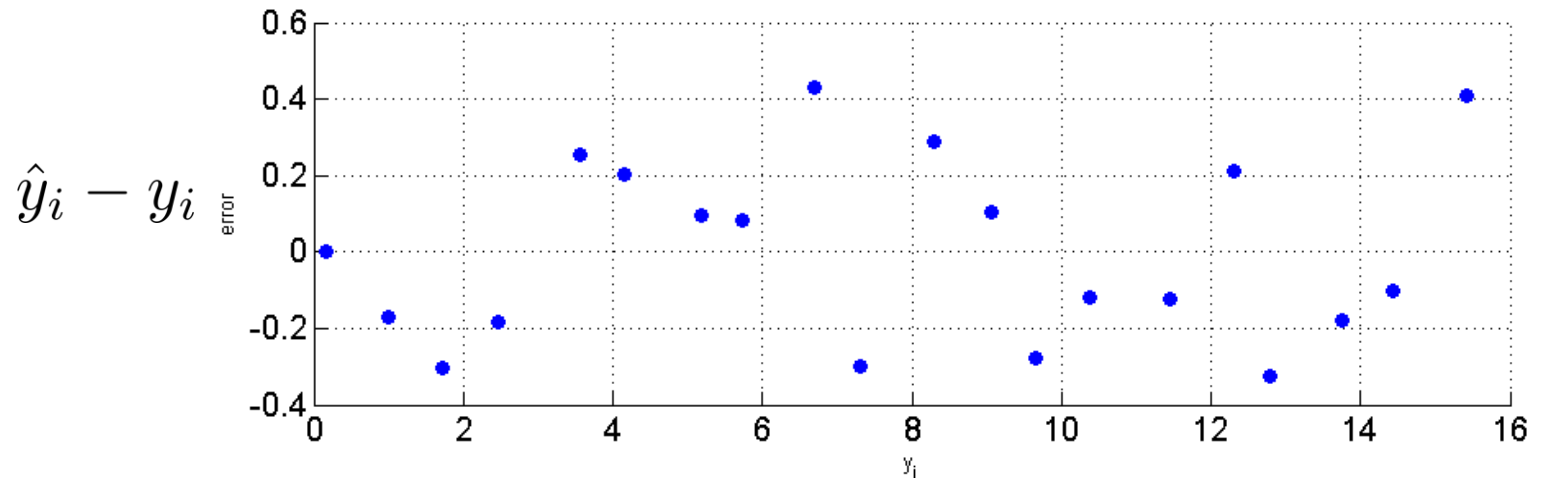




## estimated error with the outliers



## estimated error without the outliers



Bottom conforms better to  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

# Measures for Model Quality

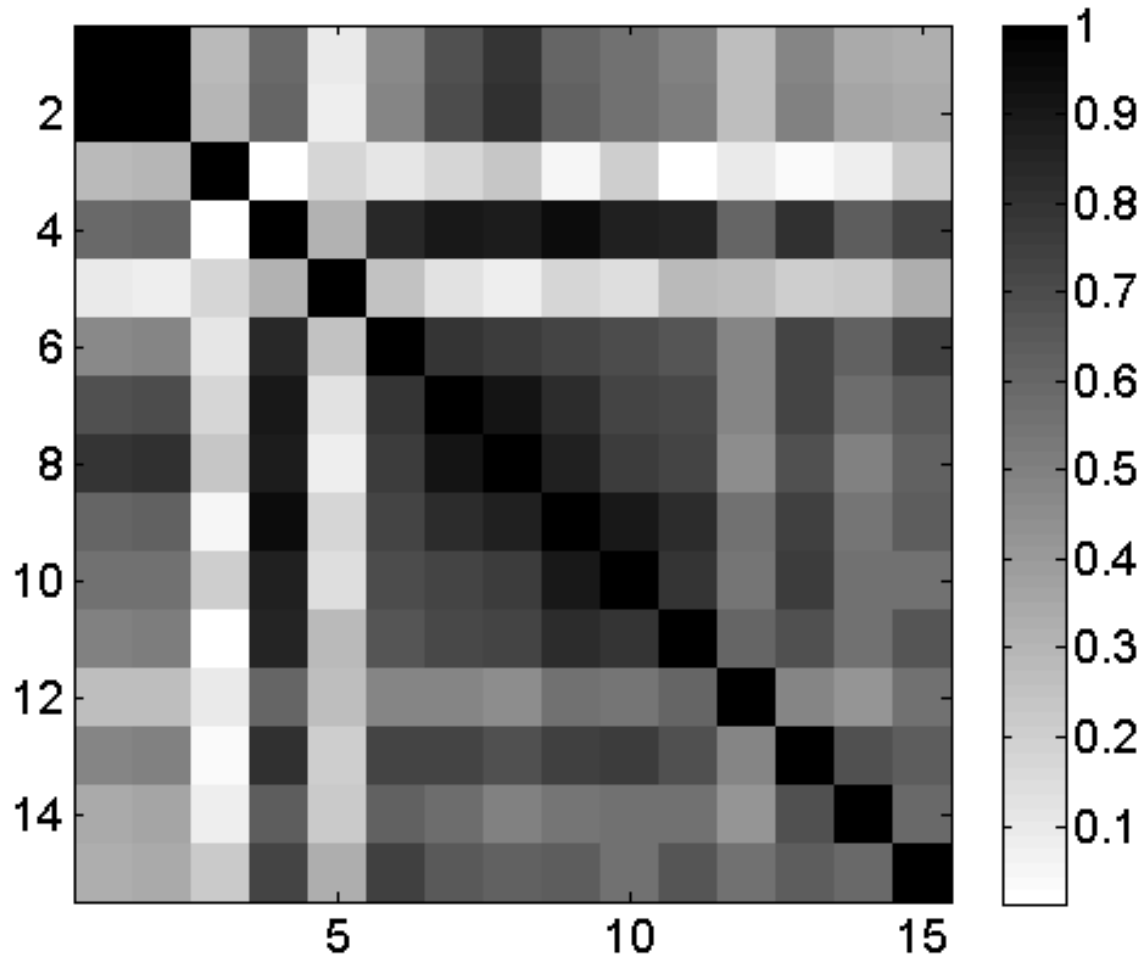
- Visual Inspection
- Mean Squared Error (MSE)  $\frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$
- Root Mean Squared Error (RMSE)  $\sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2}$
- Coefficient of determination

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{\sum_{i=1}^N (y_i - \bar{y}_i)^2}$$

- Number of Parameters
- Correlation (pairwise)
- Statistics of the Parameters:
  - ▣ Confidence Interval  $\hat{w}_i \pm \hat{\sigma}_i, \quad \hat{w}_i \pm 2\hat{\sigma}_i$
- ...

# Correlation

Example: absolute values of correlation matrix



# Diagnostic Insights

Improve model by

- Check model requirements
  - Define outliers: there is no unique definition
  - Estimate model on subset of points
  - Inspect changes of
    - ▣ Errors (e.g. studentized residuals)
    - ▣ estimated parameters
  - Consider changes of their values and statistics
- Bias vs. variance

# Multivariate Regression

What changes from one to many ?

$$y_i = w_0 + w_1 x_i + \epsilon_i, \quad i = 1, \dots, N$$

$$\rightsquigarrow y_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_D x_{iD} + \epsilon_i, \quad i = 1, \dots, N$$

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{X} := \begin{pmatrix} 1 & x_{11} & \cdots & x_{1D} \\ 1 & x_{21} & \cdots & x_{2D} \\ & \vdots & & \\ 1 & x_{N1} & \cdots & x_{ND} \end{pmatrix}, \quad \mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{pmatrix}$$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

$$\min \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \underbrace{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\mathbf{H}} \mathbf{y}$$

The same estimator as before can be used! ☺ 61

# Multivariate Regression

What is a „linear model“?

$$(1) \quad y_i = w_0 + w_1 x_i + \epsilon_i$$

$$(2) \quad y_i = w_0 + w_1 x_{1i} + w_2 x_{2i} + \epsilon_i$$

$$(3) \quad y_i = w_0 + w_1 x_i + w_2 x_i^2 + \epsilon_i$$

$$(4) \quad y_i = w_0 + w_1 x_i + w_2^2 x_i + \epsilon_i$$

$$(5) \quad y_i = w_0 + w_1 \log(x_i) + w_2 x_i + \epsilon_i$$

$$(6) \quad y_i = w_0 + w_1 \log(x_i) + w_2^3 x_i + \epsilon_i$$

$$\rightsquigarrow \mathbf{y} = \mathbf{Z}\mathbf{w} + \boldsymbol{\epsilon}$$

1, 2, 3, 5 are linear with respect to the parameters  $\mathbf{w}$

# Multivariate Regression

- Multivariate linear model as

$$y_i = w_0 + w_1 x_{i1} + w_2 x_{i2}^2 + \epsilon_i, \quad i = 1, \dots, N$$

$$y_i = w_0 + w_1 z_{i1} + w_2 z_{i2} + \epsilon_i, \quad i = 1, \dots, N$$

- Multivariate polynomial model

with new higher-order variables

(basis functions, kernel trick: Chapter 13)

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{pmatrix}, \quad \mathbf{Z} := \begin{pmatrix} 1 & x_{11} & x_{11}^2 \\ 1 & x_{21} & x_{21}^2 \\ \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N1}^2 \end{pmatrix} \in \mathbb{R}^{N \times M}$$

$$\mathbf{y} = \mathbf{Z}\mathbf{w} + \boldsymbol{\epsilon}$$

The same estimator as before can be used! ☺

# Multivariate Regression

$$y_i = w_0 + w_1 z_{i1} + w_2 z_{i2} + \dots + w_M z_{iM} + \epsilon_i, \quad i = 1, \dots, N$$

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{pmatrix}$$

$$\mathbf{Z} := \begin{pmatrix} 1 & x_{11} & x_{11}^2 & \sin(x_{11}) & \cdots & \log(x_{1P}) \\ 1 & x_{21} & x_{21}^2 & \sin(x_{21}) & \cdots & \log(x_{2P}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N1}^2 & \sin(x_{N1}) & \cdots & \log(x_{NP}) \end{pmatrix} \in \mathbb{R}^{N \times M}$$

$$\mathbf{y} = \mathbf{Z}\mathbf{w} + \boldsymbol{\epsilon}$$

The same estimator as  
before can be used! ☺

$$\hat{\mathbf{w}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{Z}\hat{\mathbf{w}} = \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$$



# Linear Models: Best Practice Tipps

- ❑ Look at the data!
- ❑ Investigate correlation
- ❑ Correlation is not causality
- ❑ Investigate errors: is there a visible trend?
- ❑ Outlier: there is no unique definition
- ❑ Point with high influence is not necessarily an outlier
- ❑ Consider different measures for model selection
  - Not only the smallest error

# Summary

- ❑ Multivariate Data
- ❑ Multivariate Normal Distribution
- ❑ Correlation vs. Dependence
- ❑ Classification with Discriminant Functions
  
- ❑ Model Quality
- ❑ Model Selection
- ❑ Preprocessing is important
- ❑ Linear Models for Regression

# Books

- *Pattern Classification*  
Duda, Hart, Stork
- *Pattern Recognition And Machine Learning*  
- Bishop
- *Regression: Models, Methods and Applications*  
- Fahrmeir, Kneib, Lang

# APPENDIX

# Covariance / Correlation

$$\mathbf{x} \in \mathbb{R}^D \quad \text{Cov}(x_i, x_j) = \sigma_{ij} \quad \text{Corr}(x_i, x_j) = \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j} \in [-1, 1]$$

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_D]^T \in \mathbb{R}^D$$

$$\begin{aligned} \boldsymbol{\Sigma} &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T - \boldsymbol{\mu}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T - 2\boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - 2\boldsymbol{\mu}\mathbb{E}[\mathbf{x}]^T + \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - 2\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T \end{aligned}$$

# Revisit: Linear Model for Regression

## Diagnostic Insights

Improve model quality

- Estimate model on subset of points
- Inspect changes of
  - ▣ Errors (studentized residuals)
  - ▣ estimated parameters
- Consider changes of their values and statistics
  - ▣ Confidence intervals  $\hat{\epsilon}_i \sim \mathcal{N}(0, \hat{\sigma}^2)$

$$\begin{aligned} X \sim \mathcal{N}(\mu, \sigma^2) &\Rightarrow P\left(\left|\frac{X - \mu}{\sigma}\right| < 2\sigma\right) = 95.5\% \\ &\Rightarrow P\left(\left|\frac{X - \mu}{\sigma}\right| < 3\sigma\right) = 99.7\% \end{aligned}$$

# Model Quality for Selection

## Confusion matrix

In general:

Positive = identified  
negative = rejected

	is positive	is negative
Predicted positive	TP	FP (Type I error)
Predicted negative	FN (Type II error)	TN

Therefore:

- True positive (TP) = correctly identified
- False positive (FP) = incorrectly identified
- True negative (TN) = correctly rejected
- False negative (FN) = incorrectly rejected