

# Linear classification

Konstantin Tretyakov (kt@ut.ee)

So far...

---

- ▶ Machine learning is important and interesting
- ▶ The general concept:

## Fitting models to data

- ▶ Some models:  
?

So far...

---

- ▶ Machine learning is important and interesting
- ▶ The general concept:

## Fitting models to data

- ▶ Some models:
  - ▶ Decision trees: if  $(x < 1)$  then  $y = 2, \dots$
  - ▶ Linear regression  $y = 2x_1 - 3x_2 + 0.5$



# Today

---

▶ **Another model:**

▶ **Linear classification**



# Linear classification

---

## ▶ Model type: **Binary\*** classification

Input: real vectors	Output: class
(1.0, -2.0, 3.1)	“ ”
(-2.1, 5.1, 3.0)	“- ”
(0.1, 3.4, -2.0)	“ ”

## ▶ Model form:

$$y = \text{sign}(2x_1 - 3x_2 + x_3 + 0.5)$$

# Linear classification

---

▶ **Binary\* classification**

# Linear classification

---

▶ **Binary\* classification**

**Can be generalized to multi-class case:**

**One-versus-All**

**All-versus-All**

# Linear classification

---

## ▶ Model type: **Binary\*** classification

Input: real vectors	Output: class
(1.0, -2.0, 3.1)	“ ”
(-2.1, 5.1, 3.0)	“- ”
(0.1, 3.4, -2.0)	“ ”

## ▶ Model form:

$$y = \text{sign}(2x_1 - 3x_2 + x_3 + 0.5)$$



# Linear classification

---

► **Model form:**

$$y = \text{sign}(w_1x_1 + w_2x_2 + w_3x_3 + w_0)$$

# Linear classification

---

## ► Model form:

$$y = \text{sign}(w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0)$$

Parameters (to be estimated)



# Linear classification

## ► Model form:

$$y = \text{sign}(w_1x_1 + w_2x_2 + w_3x_3 + w_0)$$

Parameters (to be estimated)

**How?**

# Linear classification

---

## ▶ **Parameter learning algorithms:**

- ▶ Fisher discriminant
- ▶ Least-squares (and variations)
- ▶ Perceptron (and variations)
- ▶ Logistic regression
- ▶ Support vector machines
- ▶  $l_1$ -norm SVM,
- ▶ Naïve Bayes\* ...



# Linear classification

---

## ▶ **Parameter learning algorithms:**

▶ Fisher discriminant

### **Important**

There are **many** parameter-fitting algorithms for the **same** model!

- ▶ 1-norm SVM,
- ▶ Naïve Bayes\* ...



# Practice

---

▶ **Logistic regression:** (NB:  $y \in \{0,1\}$ )

- ▶ `m = glm(y ~ X, family=binomial)`
- ▶ `predict(m, newX)`

▶ **SVM:** (NB:  $y \in \{-1,1\}$ )

- ▶ `library('e1071')`
- ▶ `m = svm(X, y, kernel='linear')`
- ▶ `predict(m, newX)`



# Practice

---

## ▶ Fisher discriminant:

- ▶ `library('MASS')`
- ▶ `m = lda(X, factor(y))`
- ▶ `predict(m, newX)$x`







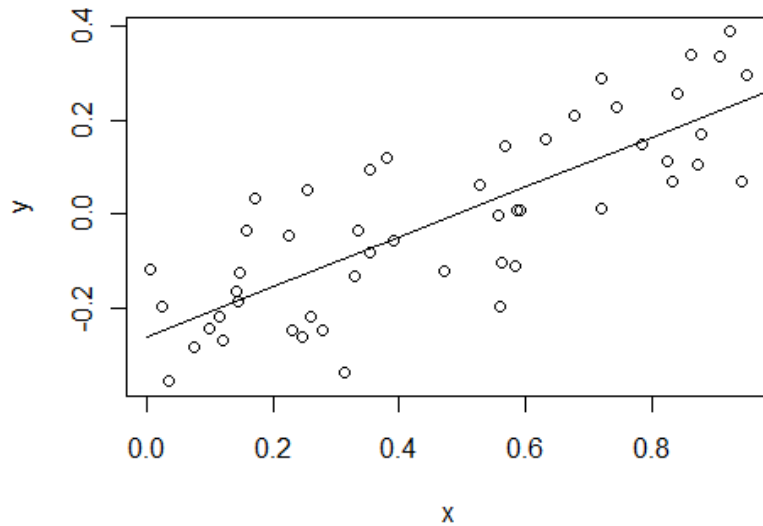
# The Boring Theory

---

- ▶ Algebra & Geometry
- ▶ Fisher's Discriminant
- ▶ Least-squares approach
- ▶ Perceptron
- ▶ Other Methods

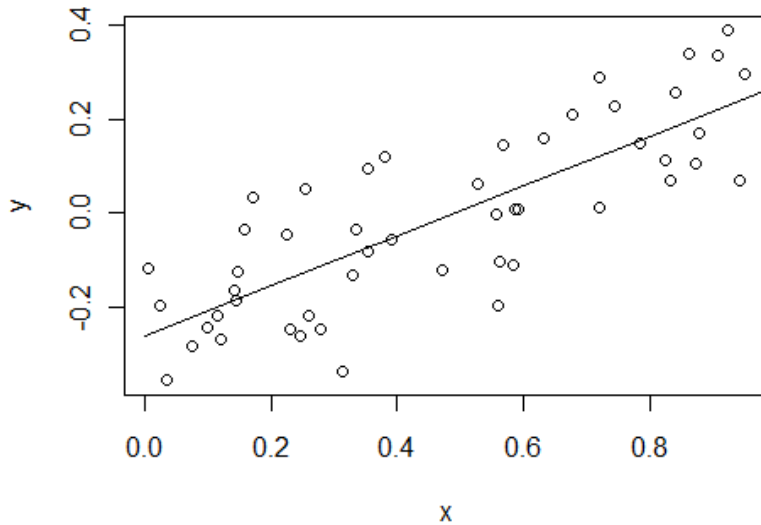
# One-dimensional case

---

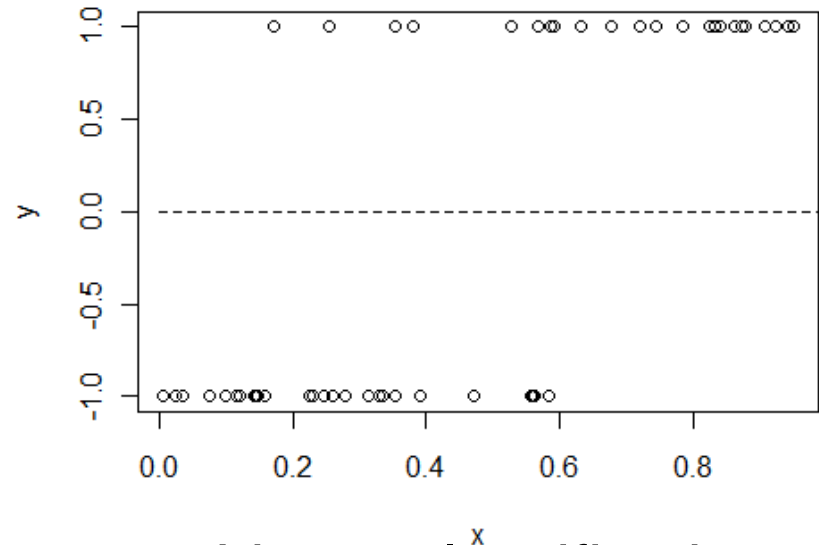


Linear regression

# One-dimensional case

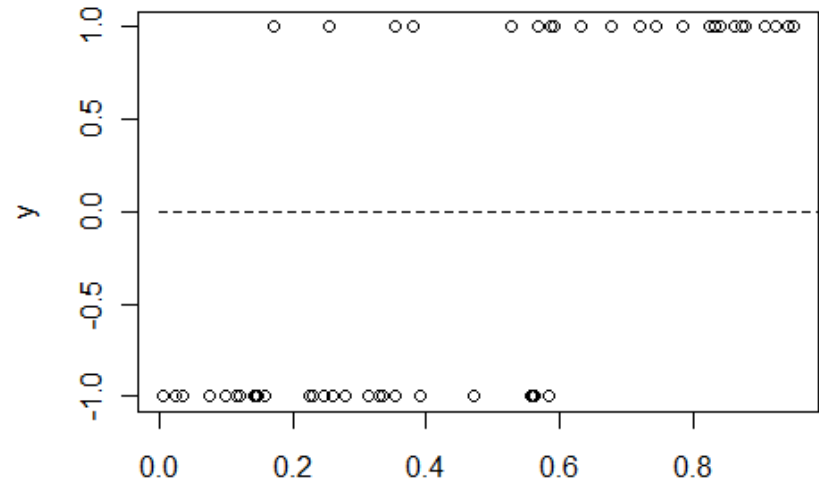
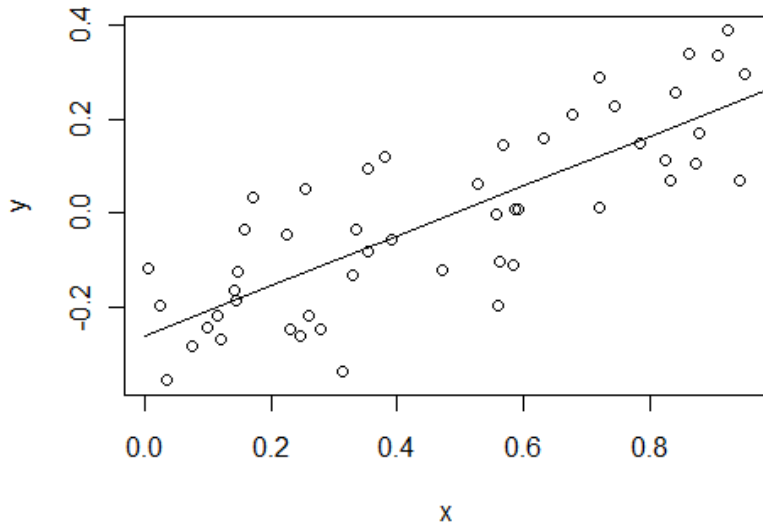


Linear regression

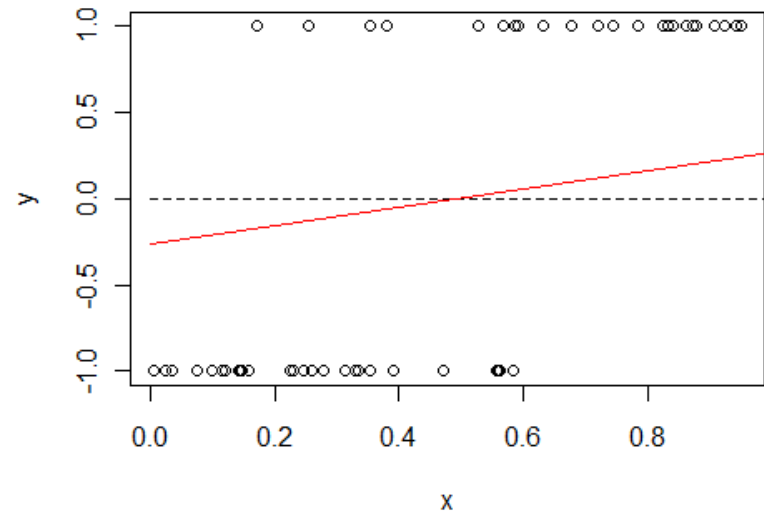


Linear classification

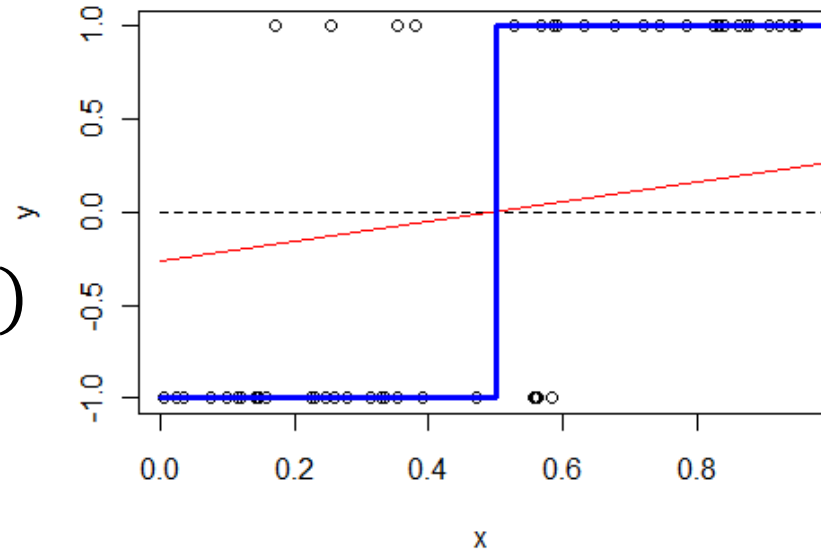
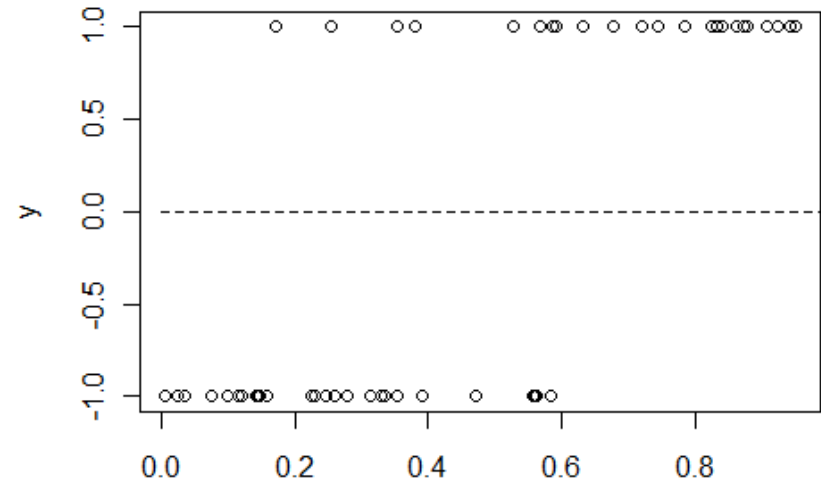
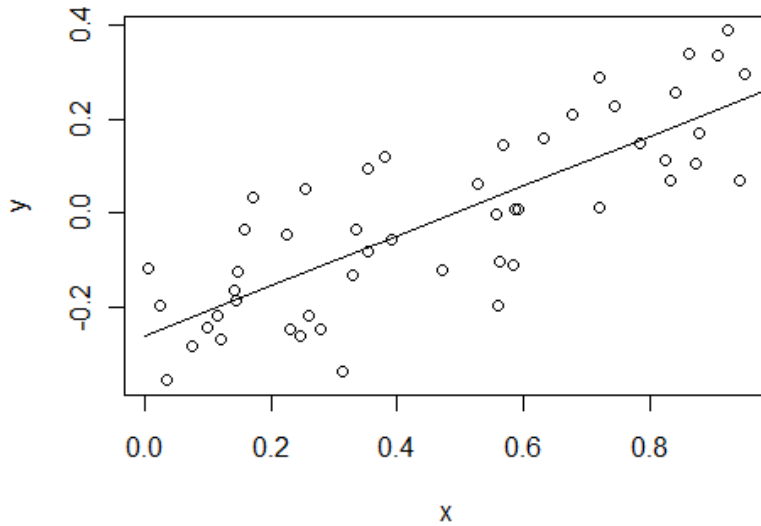
# One-dimensional case



$$y = w_0 + w_1 x$$



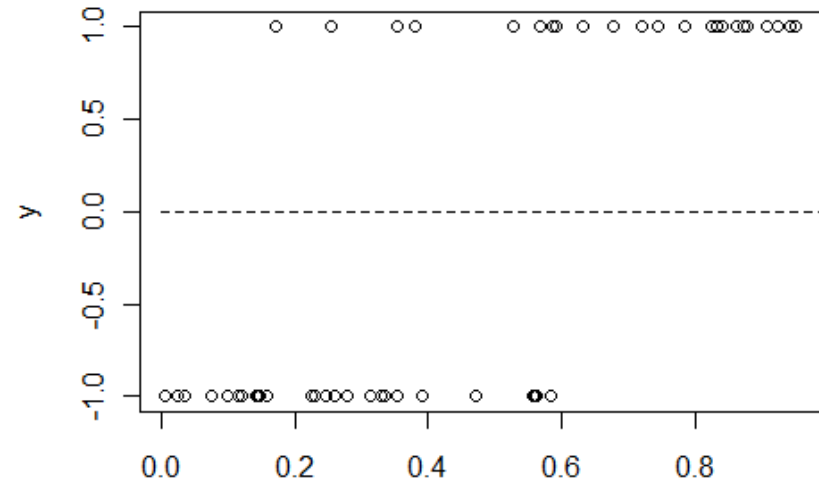
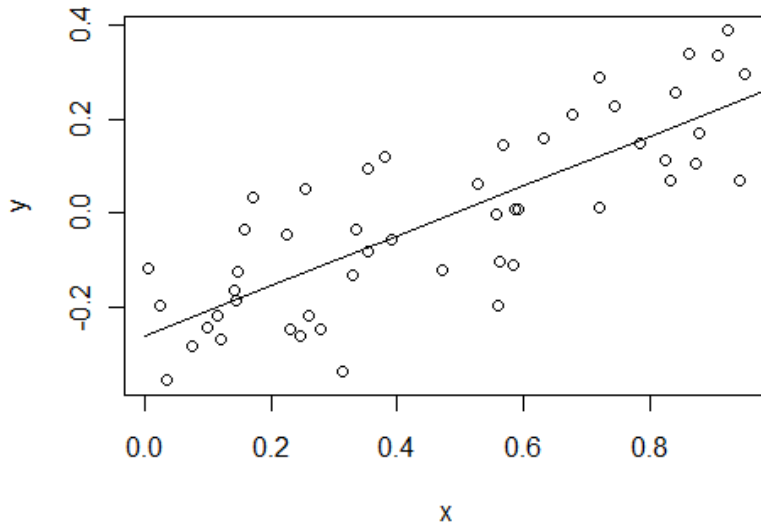
# One-dimensional case



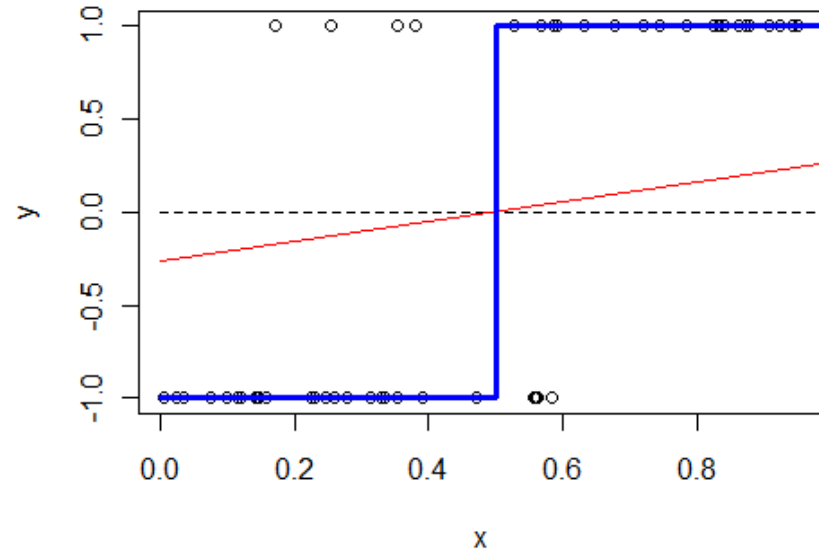
$$y = \text{sign}(w_0 + w_1 x)$$



# One-dimensional case

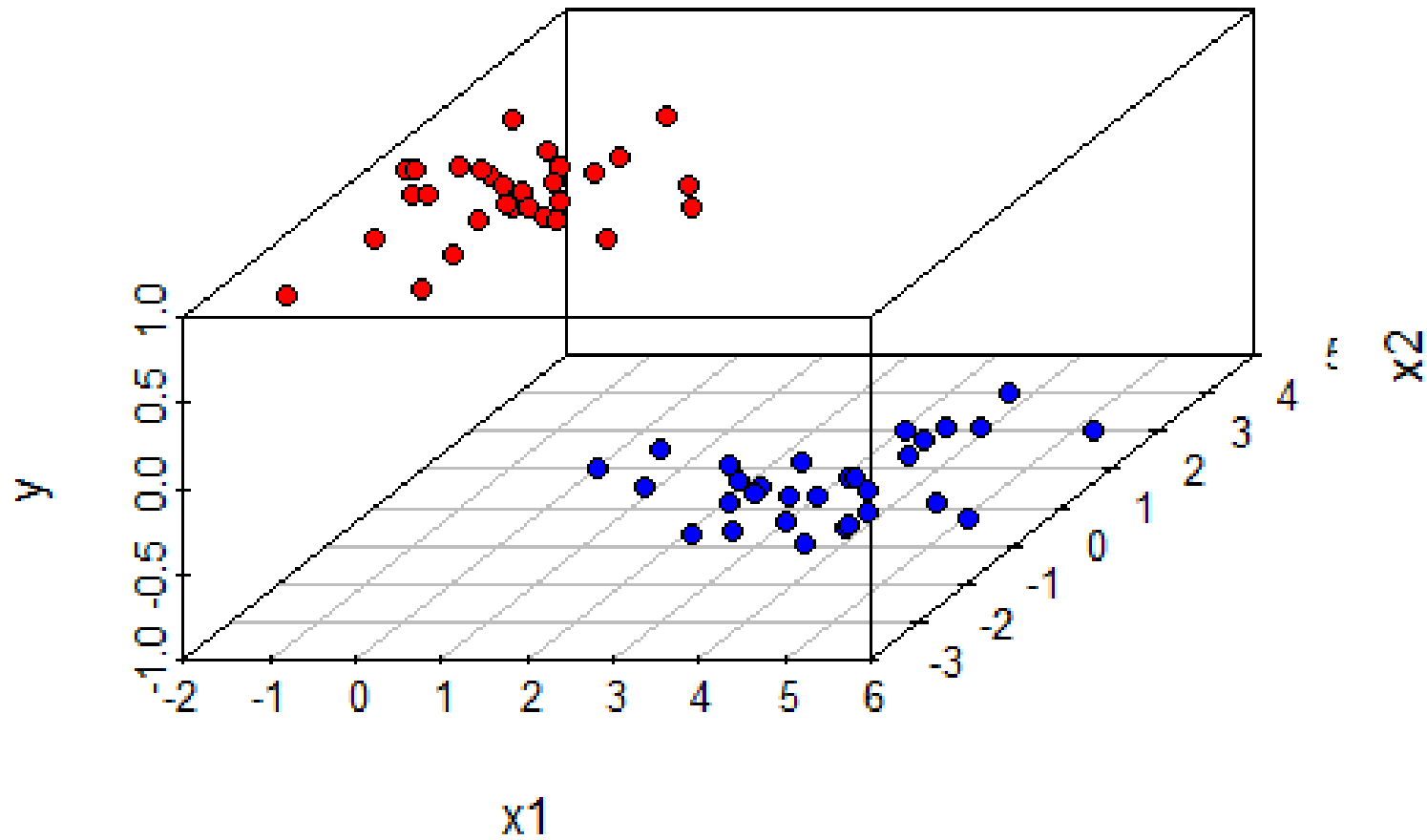


$$y = \begin{cases} 1 & \text{if } x \geq t \\ -1 & \text{otherwise} \end{cases}$$

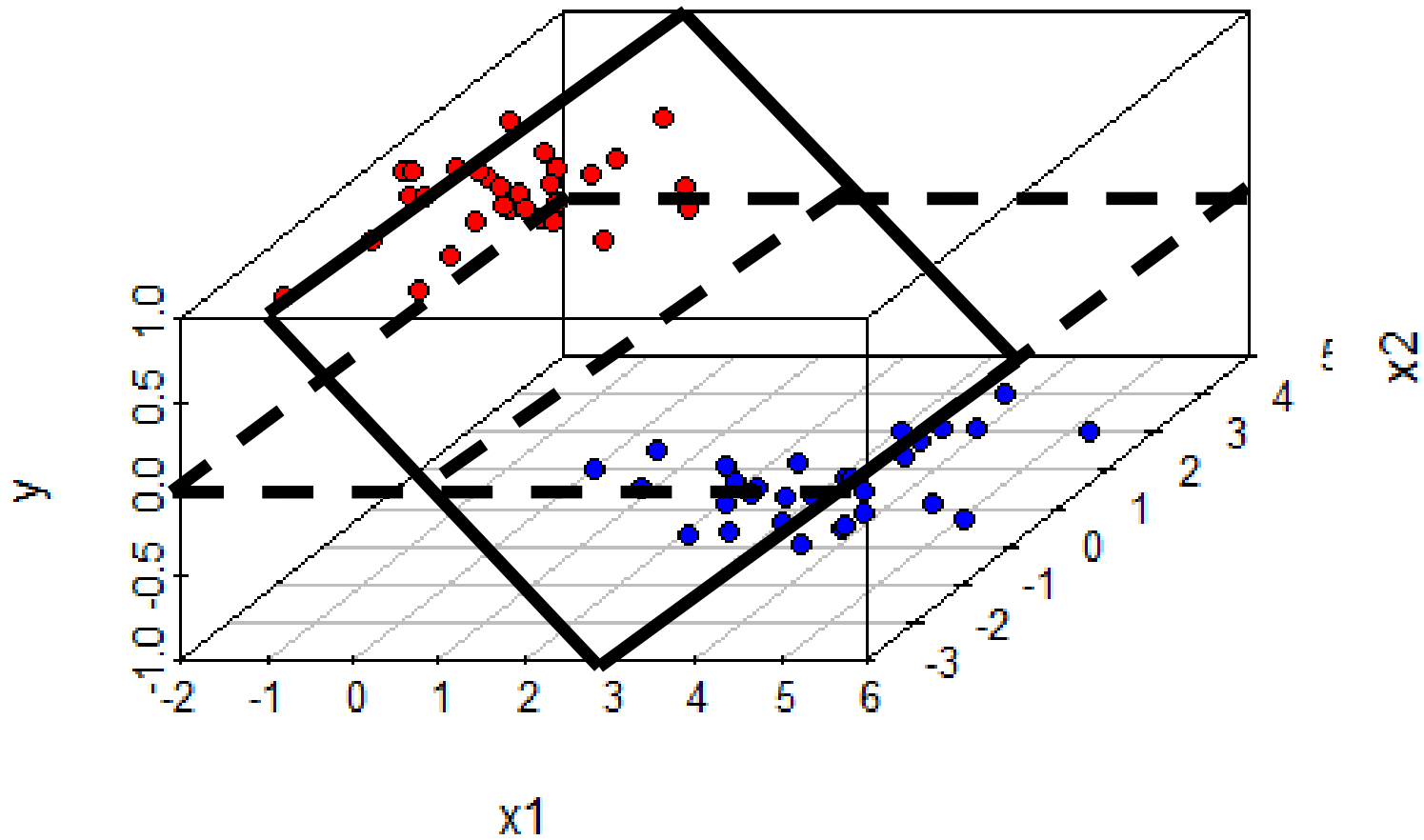


Linear classifier in a one-dimensional case is just a cutoff value.

# Two-dimensional case



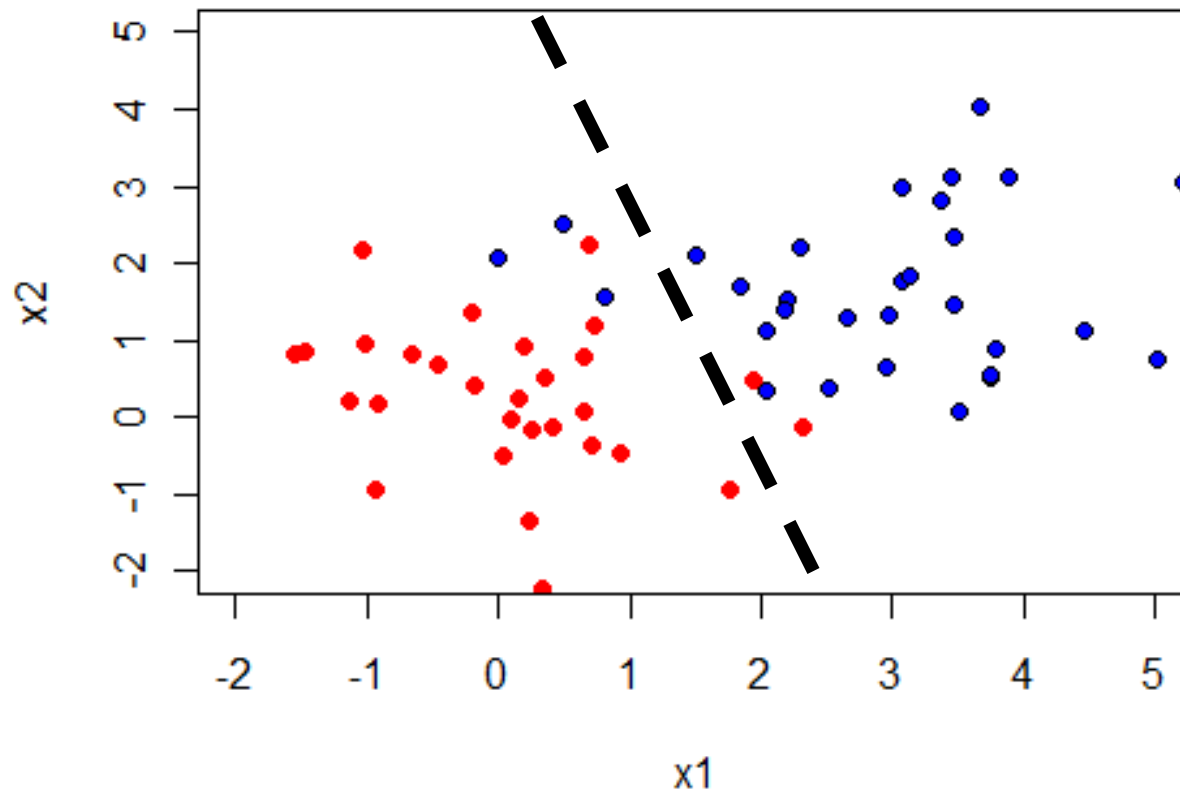
# Two-dimensional case





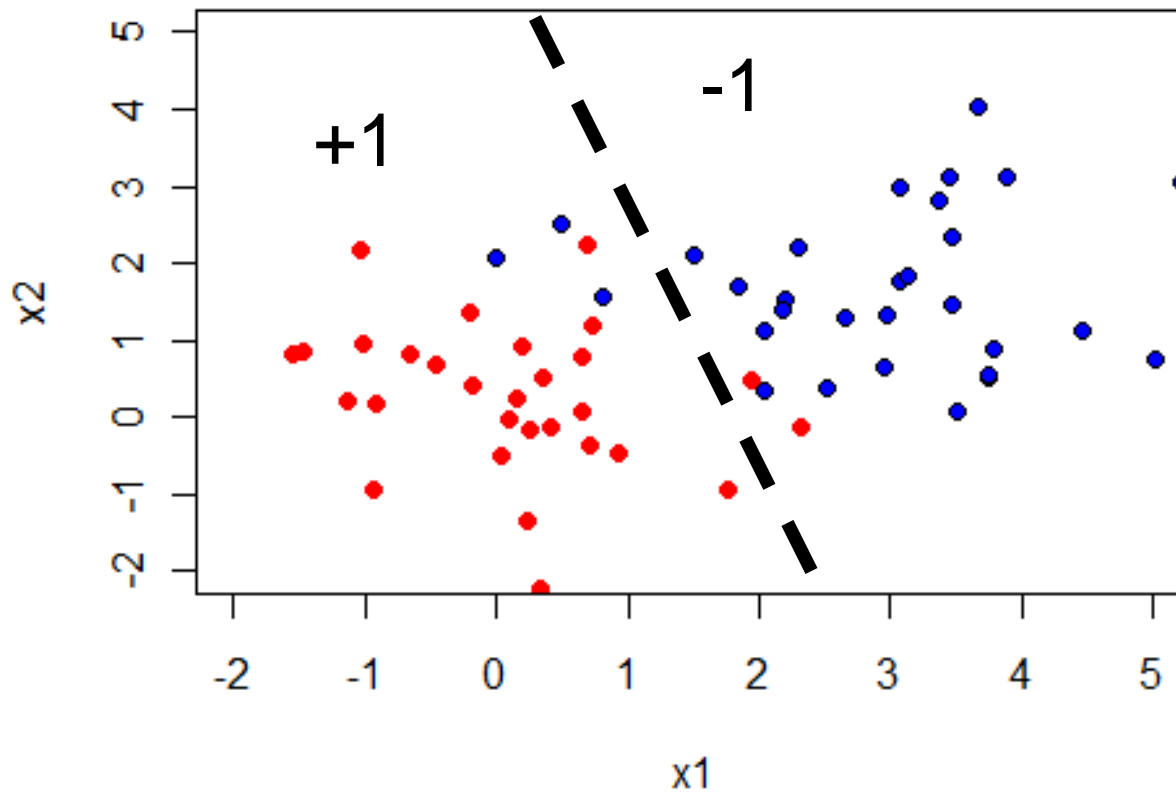
# Two-dimensional case

---



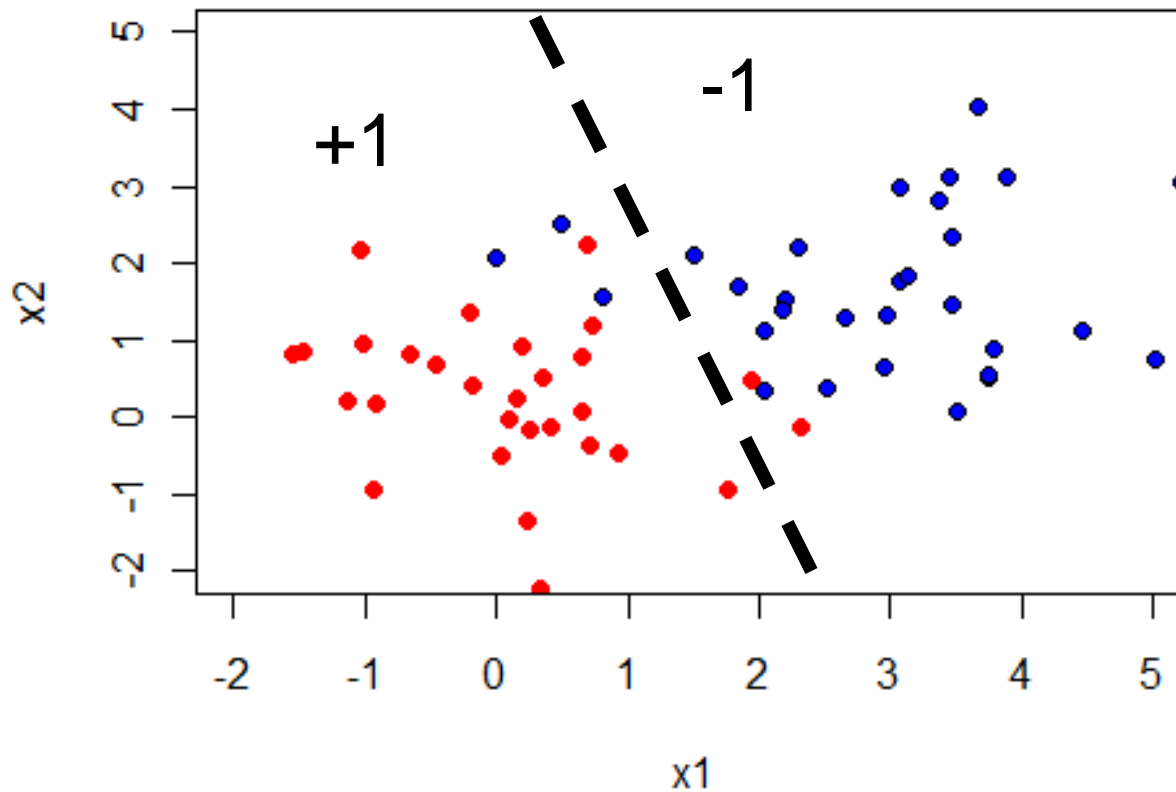
# Two-dimensional case

$$f(x_1, x_2) = \text{sign}(w_1x_1 + w_2x_2 + w_0)$$



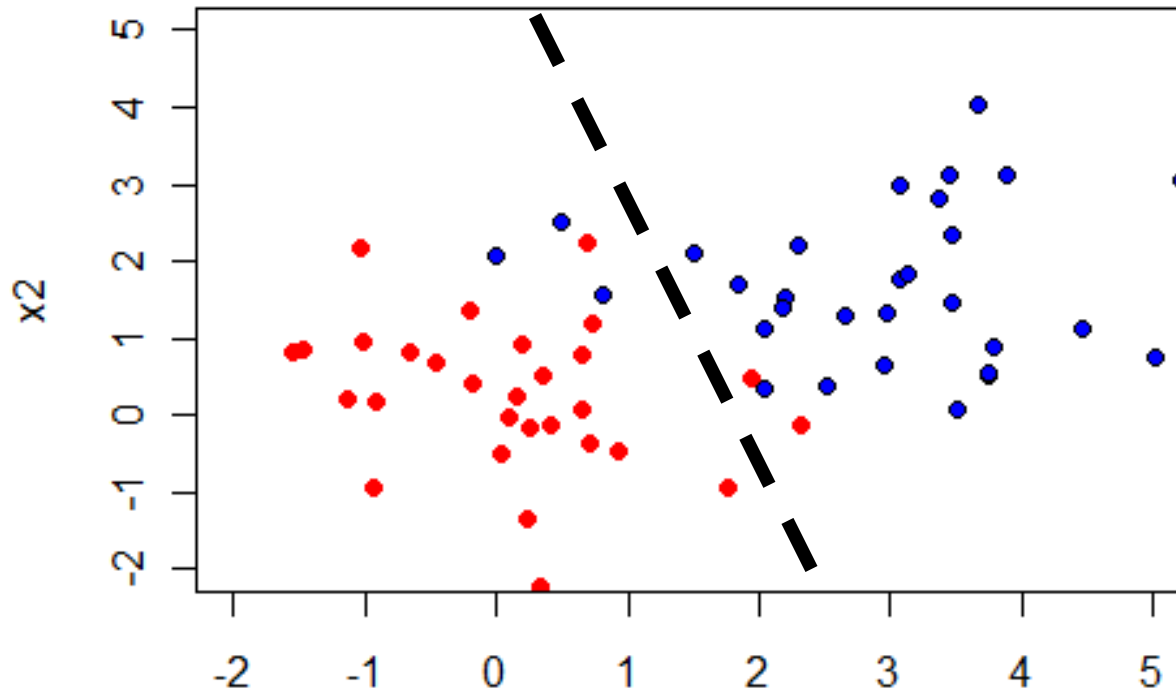
# Two-dimensional case

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0)$$

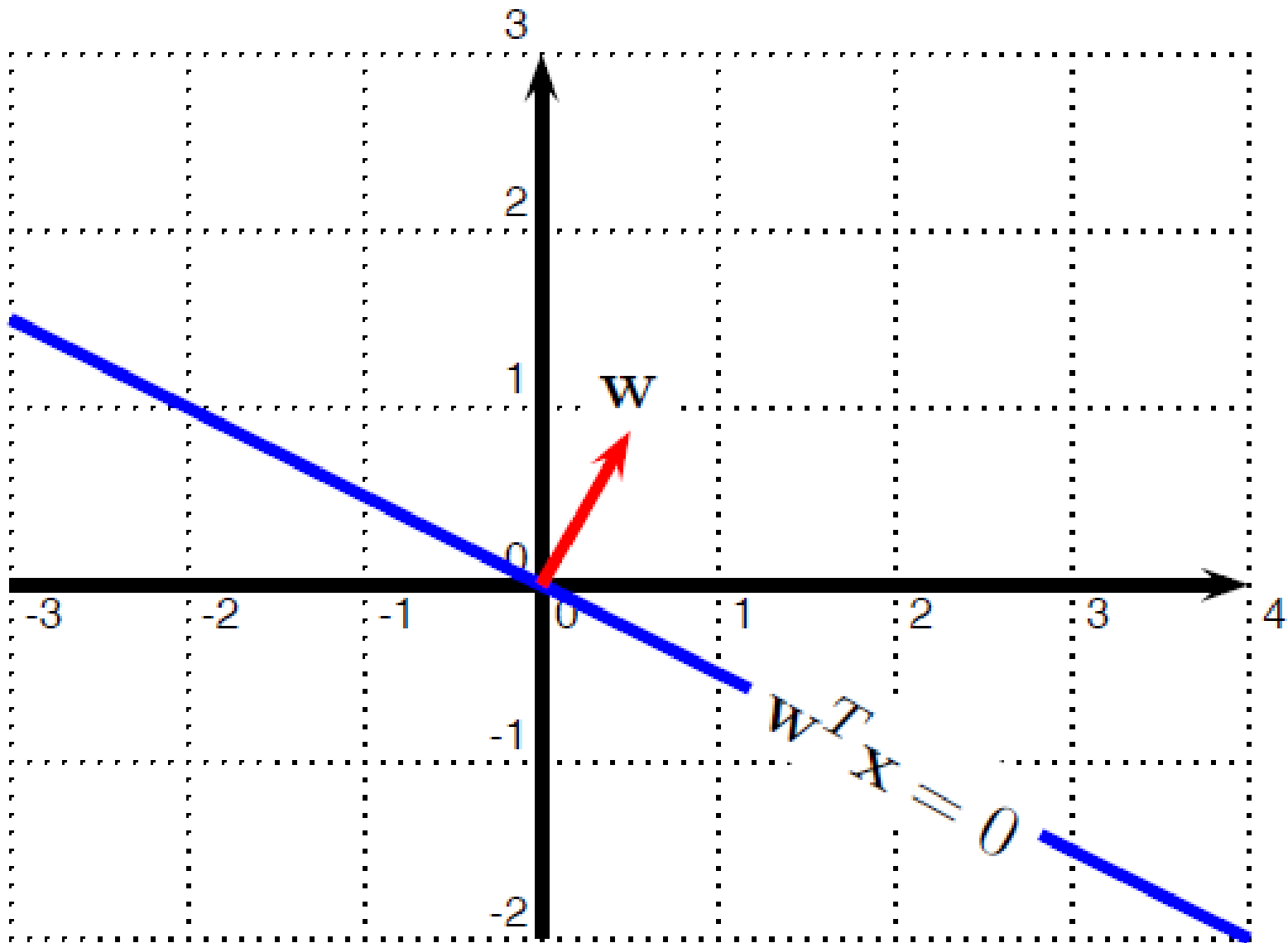


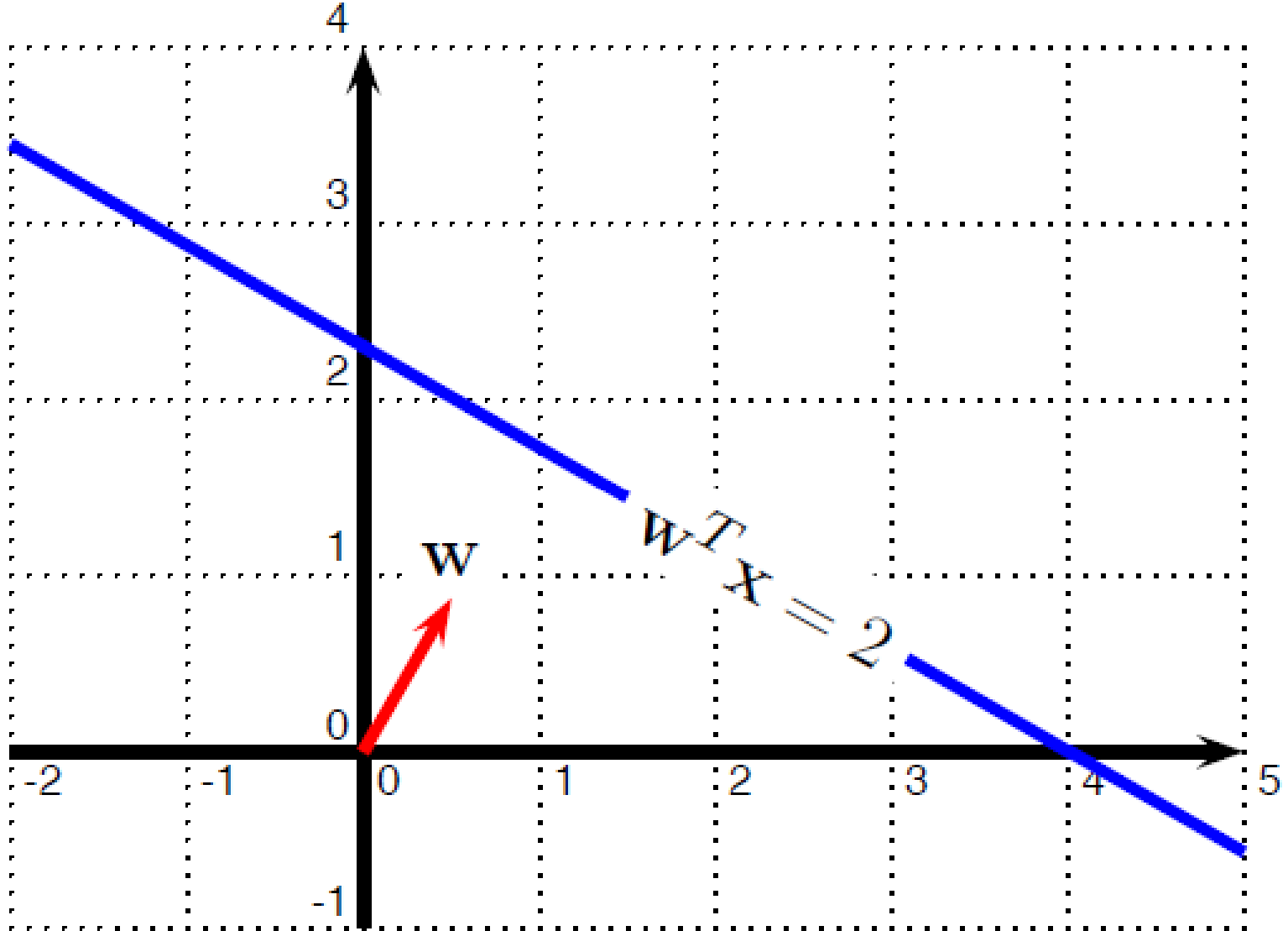
# Two-dimensional case

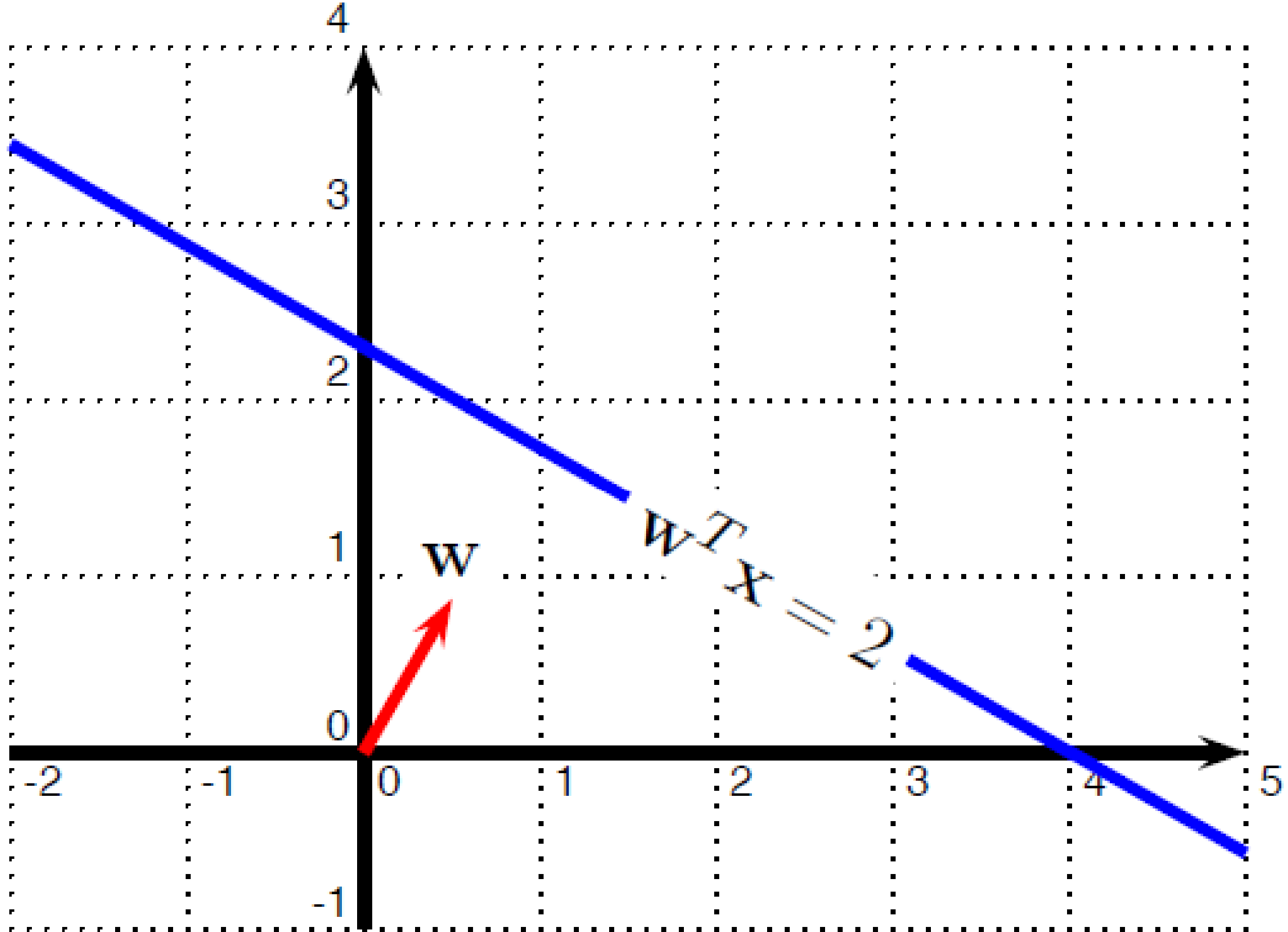
$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$



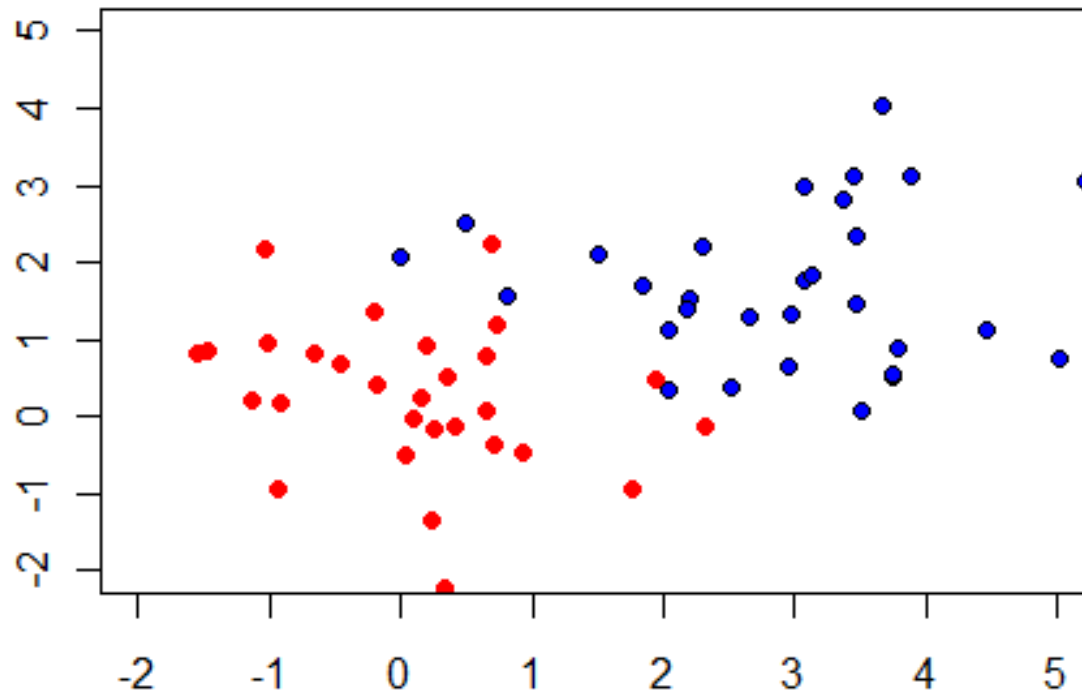
Separating line/plane/hyperplane







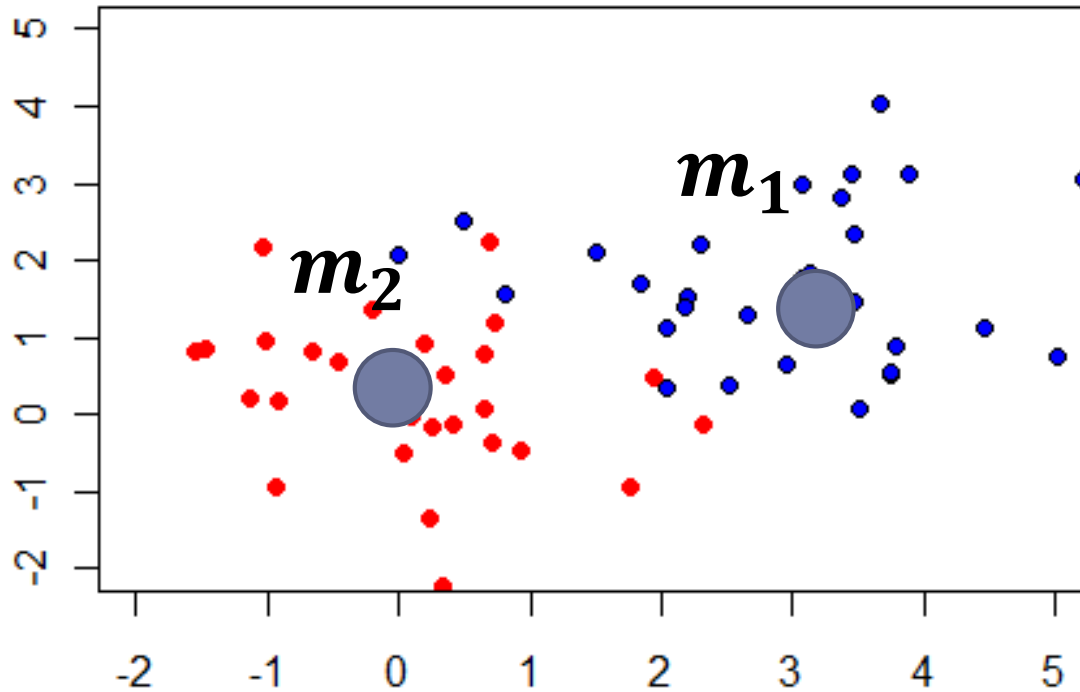
$$w^T (x - p) = 0$$





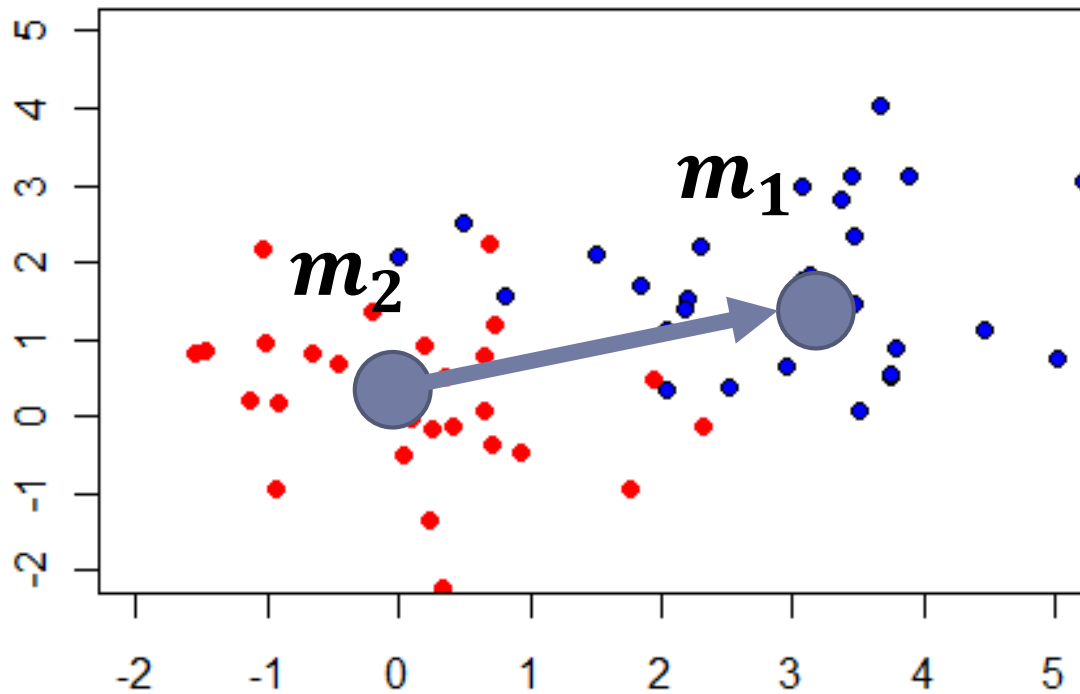
# Fisher's discriminant

---



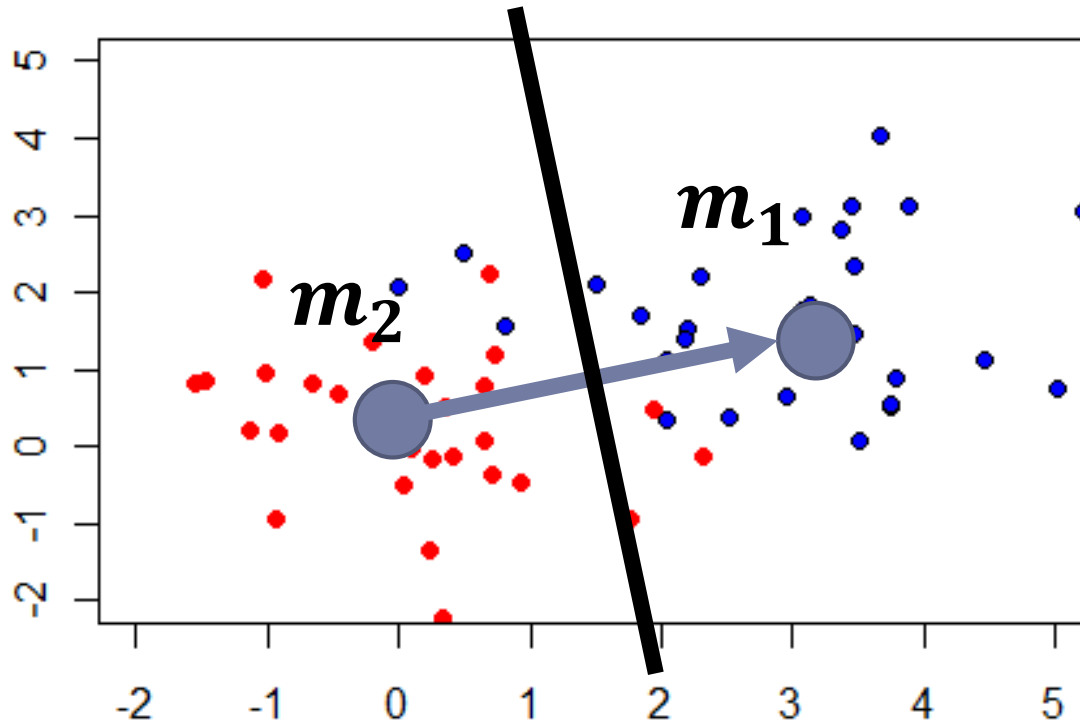
# Fisher's discriminant

---



$$w = m_1 - m_2$$

# Fisher's discriminant

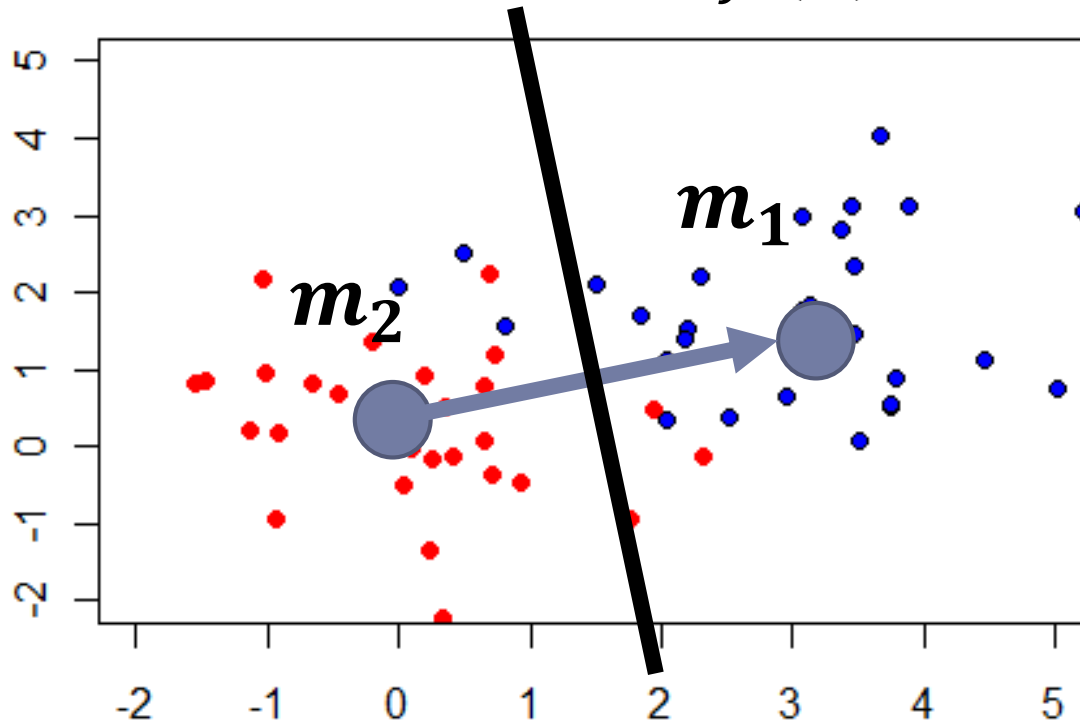


$$w = m_1 - m_2$$

$$p = (m_1 + m_2)/2$$

# Fisher's discriminant

$$f(x) = w^T (x - p)$$



$$w = m_1 - m_2$$

$$p = (m_1 + m_2) / 2$$

# Fisher's discriminant

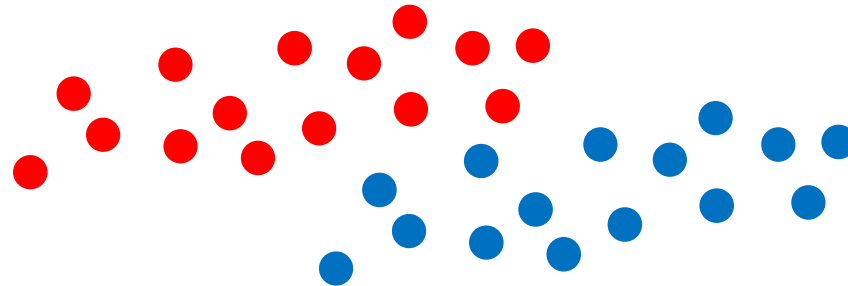
---

- ▶ Any problems with this approach?

# Fisher's discriminant

---

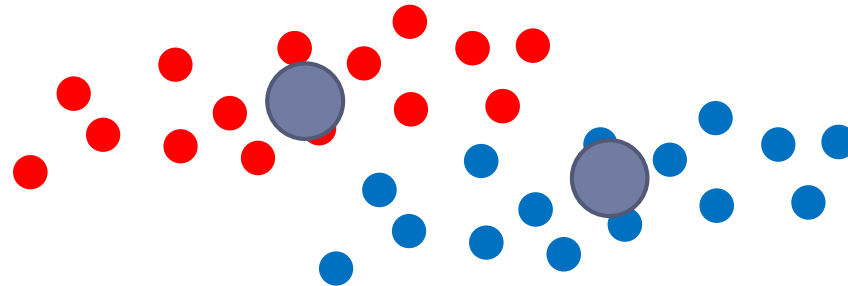
- ▶ Any problems with this approach?



# Fisher's discriminant

---

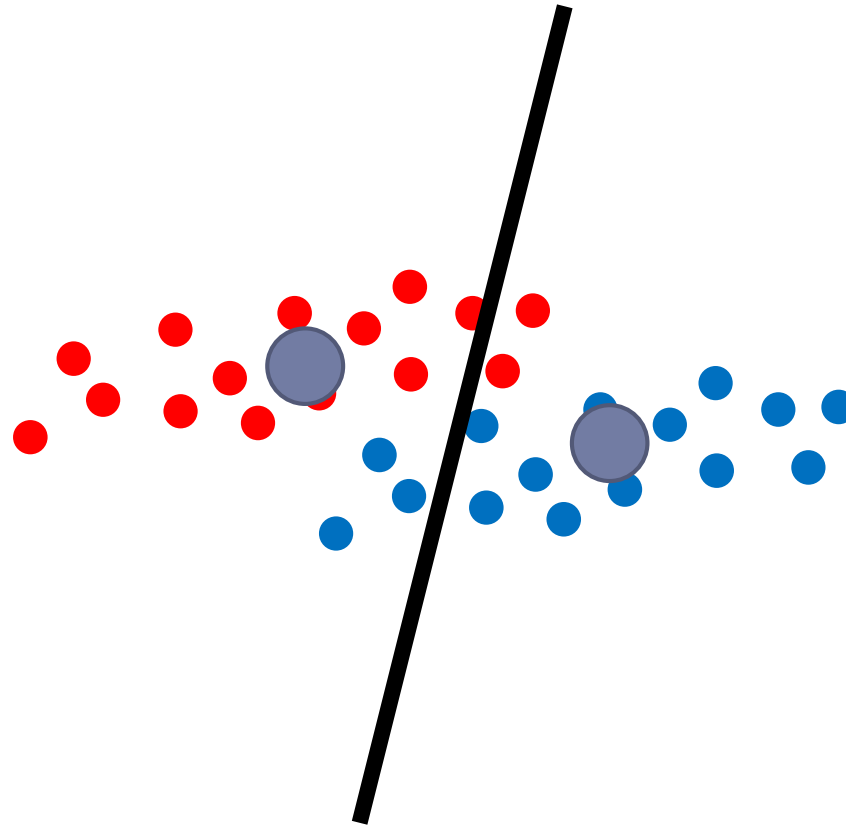
- ▶ Any problems with this approach?



# Fisher's discriminant

---

- ▶ Any problems with this approach?

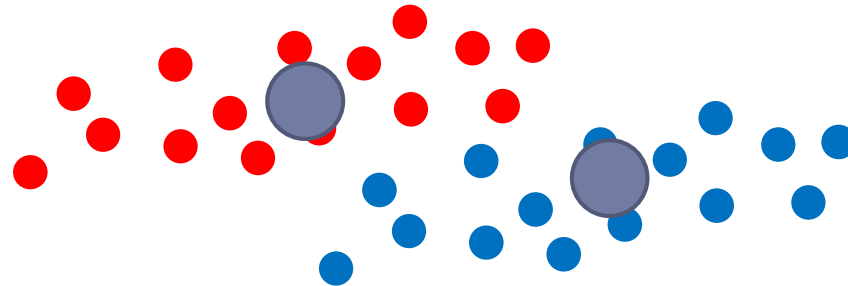




# Fisher's discriminant

---

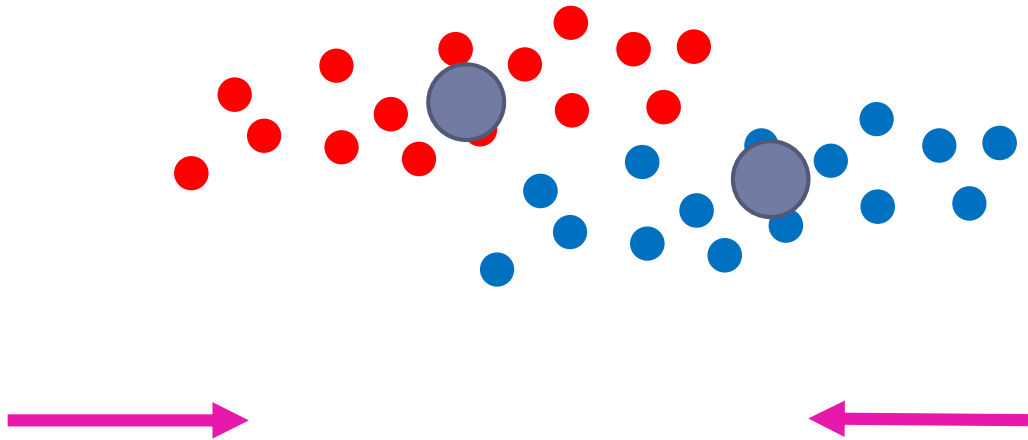
- ▶ Any problems with this approach?



# Fisher's discriminant

---

- ▶ Any problems with this approach?

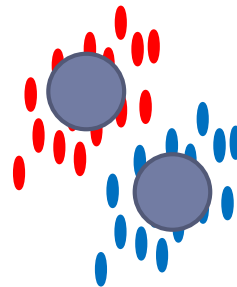


Therefore, we augment the idea by first **transforming the space** to ensure **maximally “circular” shape** of the two point clouds.

# Fisher's discriminant

---

- ▶ Any problems with this approach?

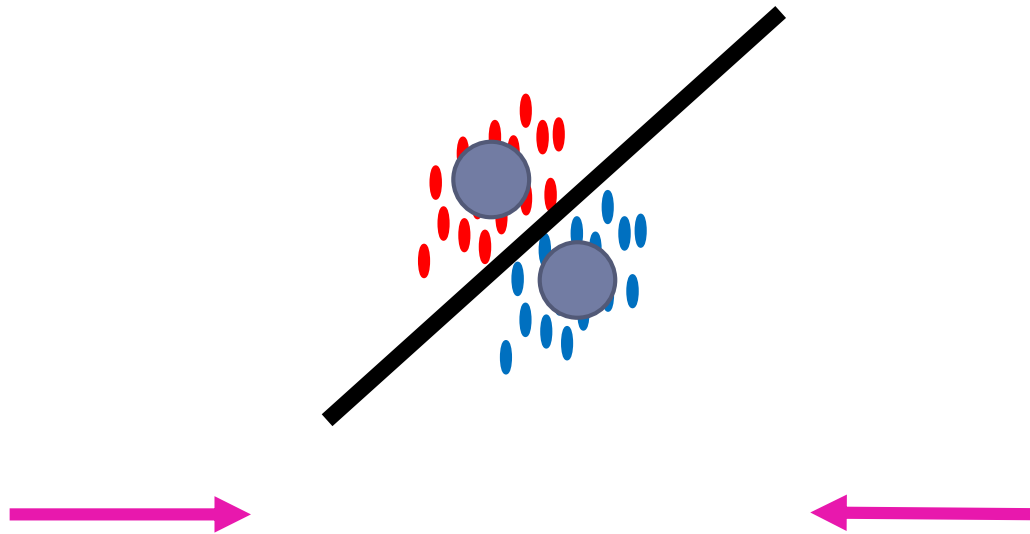


Therefore, we augment the idea by first **transforming the space** to ensure **maximally “circular” shape** of the two point clouds.

# Fisher's discriminant

---

- ▶ Any problems with this approach?



Therefore, we augment the idea by first **transforming the space** to ensure **maximally “circular” shape** of the two point clouds.

# Fisher's discriminant

---

After doing some math, the estimation rule turns out to be the following:

(we shall discuss the derivation on the practice session)

$$\mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2).$$

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2. \quad \mathbf{S}_i = \sum_{\mathbf{x} \in \mathcal{D}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t$$

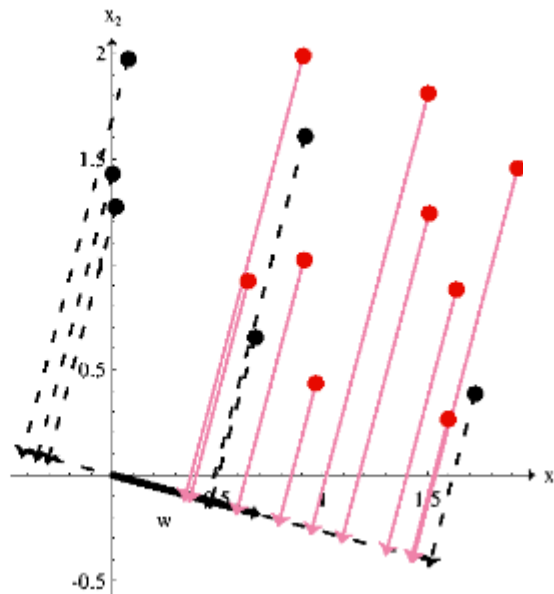
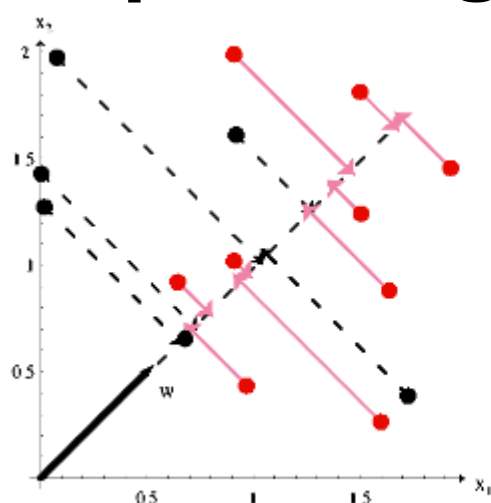
$$\mathbf{w}_0 = -\mathbf{w}^T (\mathbf{m}_1 + \mathbf{m}_2) / 2$$

# Other Interpretations of the FD

## ► Probabilistic interpretation

- Asymptotically optimal classifier for two normally-distributed classes with equal covariance.

## ► Optimizing the Fisher's criterion



$$J(\mathbf{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

# The Boring Theory

---

→ ~~Algebra & Geometry~~

→ ~~Fisher's Discriminant~~

$$\mathbf{w} = \mathbf{S}_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2).$$

▶ Least-squares approach

▶ Perceptron

▶ Other Methods



# The “bias term” convention

---

$$f(\mathbf{x}) = w_0 + \mathbf{w}^T \mathbf{x}$$

is equivalent to

$$f(\mathbf{x}') = \mathbf{w}'^T \mathbf{x}'$$

where  $\mathbf{x}'$  and  $\mathbf{w}'$  are *augmented* vectors:

$$\mathbf{x}' = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} \quad \mathbf{w}' = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$



# The “bias term” convention

---

$$f(\mathbf{x}) = w_0 + \mathbf{w}^T \mathbf{x}$$

is equivalent to

**NB: Many algorithms (e.g. Fisher’s discriminant) treat the bias term in a special way!**

where  $\mathbf{x}'$  and  $\mathbf{w}'$  are *augmented* vectors:

$$\mathbf{x}' = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} \quad \mathbf{w}' = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} =$$

# LSE Methods

---

- ▶ Recall linear regression:

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

# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

Why not apply the same idea directly to solve the classification tasks?

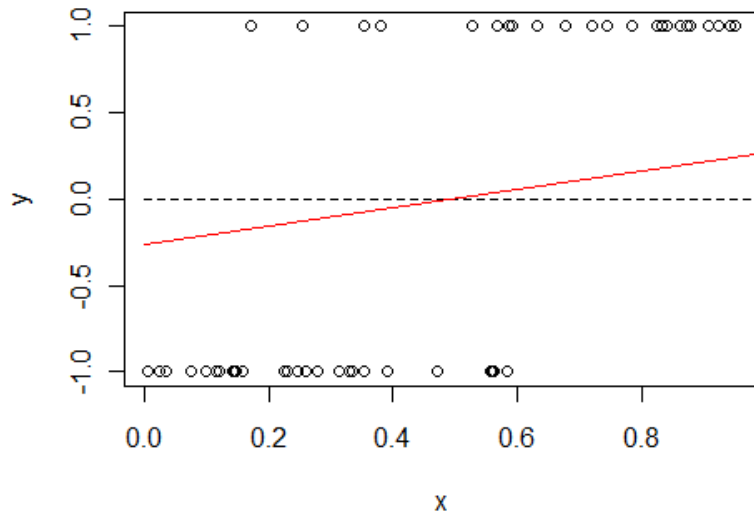
# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

Why not apply the same idea directly to solve the classification tasks?



**It makes no sense!**

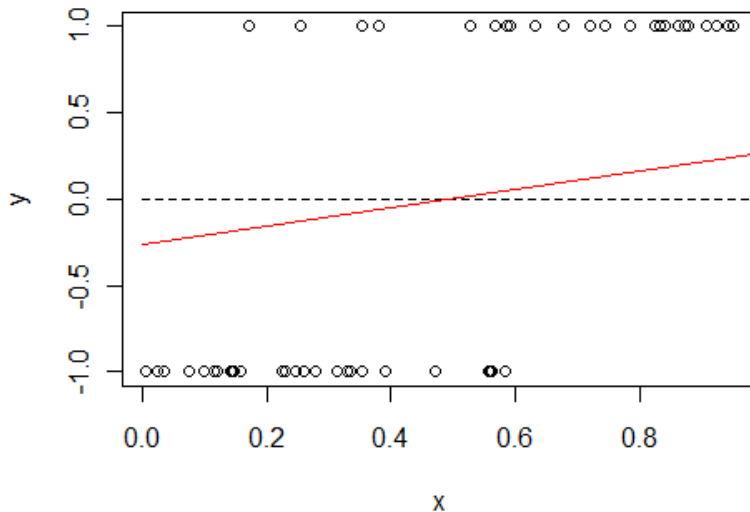
# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

Why not apply the same idea directly to solve the classification tasks?



**It is too sensitive to outliers**

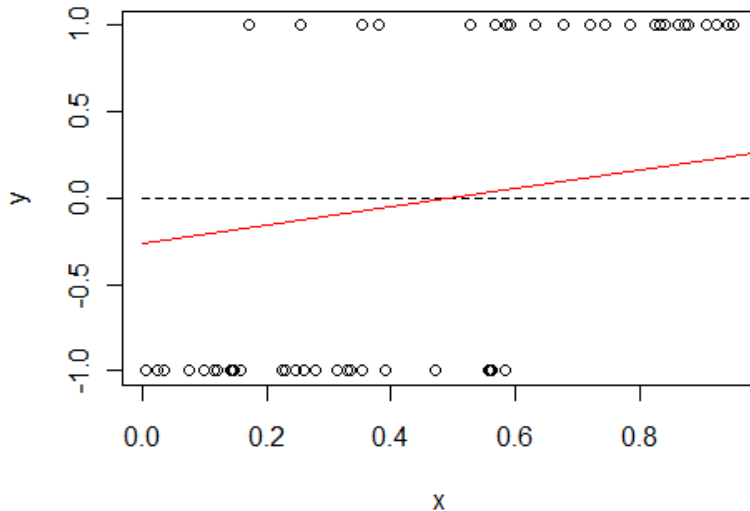
# LSE Methods

---

- ▶ Recall linear regression:

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

Why not apply the same idea directly to solve the classification tasks?



**... but it works!\***

\*sometimes



# LSE Methods – Why?

---

Suppose there exists some good classifier

$$y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i),$$

# LSE Methods – Why?

---

Suppose there exists some good classifier

$$y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i),$$

If we knew the *margin* that the classifier has at each point:

$$y'_i = \mathbf{w}^T \mathbf{x}_i$$

we could recover  $\mathbf{w}$  by fitting a linear regression on  $(y'_i, \mathbf{x}_i)$ .

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}'$$

# LSE Methods – Why?

---

Suppose there exists some good classifier

$$y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i),$$

Unfortunately, we do not know the margins.

# LSE Methods – Examples

---

- ▶ It is possible to show that if we pick margins  $y_i'$  as follows:

$$y'_i = y_i / P(y_i)$$

where  $P(c)$  = proportion of elements of class  $c$  in the sample,

then LSE result **is equivalent to the Fisher's discriminant.**

# Ho-Kashyap Algorithm

---

- ▶ It is possible to iteratively search for the margin vector, optimizing the error

$$\|X\mathbf{w} - \mathbf{y}'\|^2$$

with the condition that components of  $\mathbf{y}'$  are *never decreased* in absolute value.

$$\mathbf{y}'_0 := \mathit{rand} \cdot \mathbf{y}$$

$$\Delta \mathbf{y}'_i := \mu((X\mathbf{w} - \mathbf{y}'_i) \cdot \mathbf{y})^+ \cdot \mathbf{y}$$



# The Boring Theory

---

▶ ~~Algebra & Geometry~~

▶ ~~Fisher's Discriminant~~

$$\mathbf{w} = \mathbf{S}_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2).$$

▶ ~~Least-squares approach~~

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

▶ Perceptron

▶ Other Methods



# Perceptron

---

- ▶ Consider the error function:

$$\sum -y_i \mathbf{w}^T \mathbf{x}_i$$

where the sum is taken **over all misclassified examples**.

- ▶ Devise a batch and stochastic gradient-descent optimization procedures.



# Perceptron

---

- ▶ Consider the error function:

$$\sum -y_i \mathbf{w}^T \mathbf{x}_i$$

where the sum is taken **over all misclassified examples.**

- ▶ **Batch perceptron**

$$\Delta \mathbf{w} = \mu \sum y_i \mathbf{x}_i$$

(sum over misclassified examples)



# Perceptron

---

- ▶ Consider the error function:

$$\sum -y_i \mathbf{w}^T \mathbf{x}_i$$

where the sum is taken **over all misclassified examples**.

- ▶ **Stochastic (classic) perceptron**

$$\Delta \mathbf{w} = \mu y_i \mathbf{x}_i$$

(for randomly picked misclassified example)

# Perceptron properties

---

- ▶ Always converges for linearly-separable data (independently of the step-size!)
- ▶ Never converges for non-separable data.
- ▶  $\mu$  does not matter, but smart choices are possible. E.g. *relaxation to margin*:

$$\mu = \frac{m - \mathbf{w}^T \mathbf{x}_i y_i}{\|\mathbf{x}_i\|^2}$$

# The Boring Theory

---

→ Algebra & Geometry

→ Fisher's Discriminant

$$\mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2).$$

→ Least-squares approach

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

→ Perceptron

$$\Delta \mathbf{w} = \mu y_i \mathbf{x}_i$$

▶ Other Methods

---



# Other methods

---

- ▶ Logistic regression:

$$\sum_i y_i \log p_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - p_{\mathbf{w}}(\mathbf{x}_i))$$

- ▶ SVM:

$$\sum_i \text{hinge\_loss}(y_i, \mathbf{w}_i^T \mathbf{x}_i) + \lambda \|\mathbf{w}_i\|^2$$

- ▶ ...

# The Boring Theory

---

→ Algebra & Geometry

→ Fisher's Discriminant

$$\mathbf{w} = \mathbf{S}_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2).$$

→ Least-squares approach

$$\mathbf{w} = \mathbf{X}^+ \mathbf{y}$$

→ Perceptron

$$\Delta \mathbf{w} = \mu y_i \mathbf{x}_i$$

→ Other Methods

other objective funcs



# Quiz

---

- ▶ How many linear classification algorithms were mentioned in the lecture?
- ▶ How many of them could you implement?
- ▶ How many of them could you use from R?

