

Linear classification

Konstantin Tretyakov (kt@ut.ee)

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**

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

▶ **Parameter learning algorithms:**

- ▶ Fisher discriminant
- ▶ Least-squares (and variations)
- ▶ Perceptron (and variations)
- ▶ Logistic regression
- ▶ Support vector machines
- ▶ l_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`



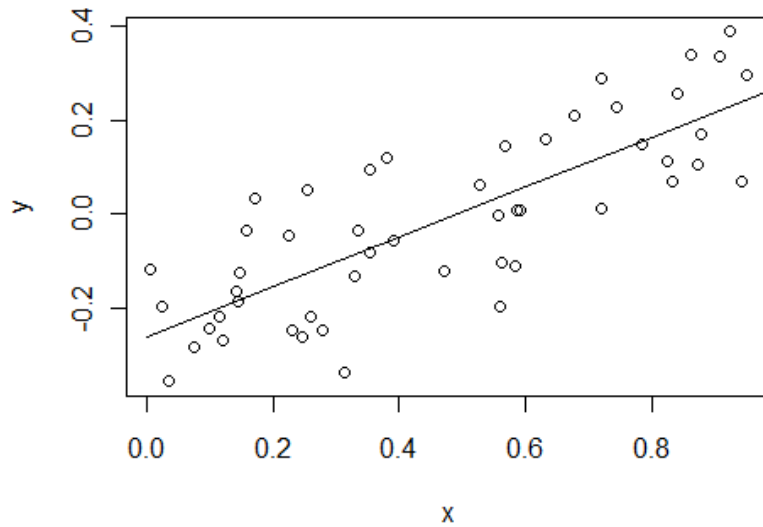


March 6, 2011

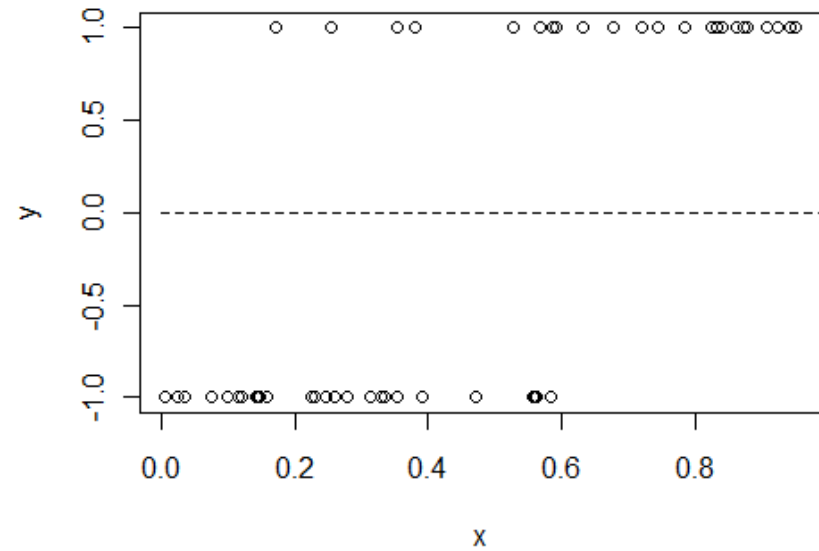
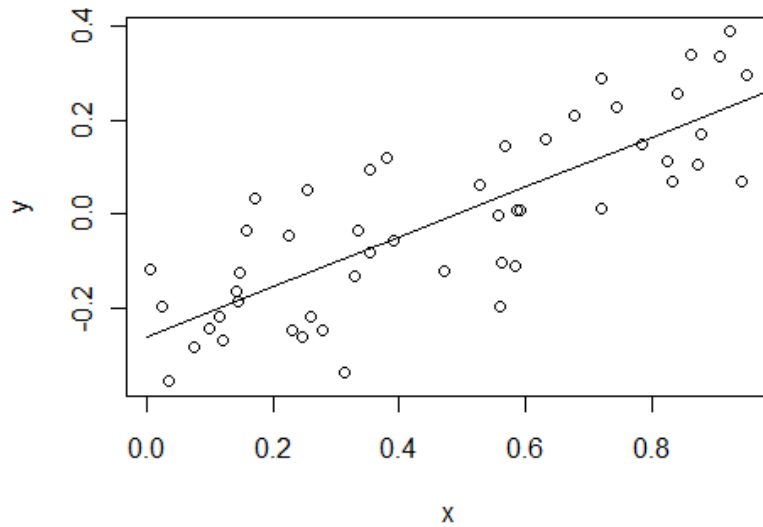
The Boring Theory

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

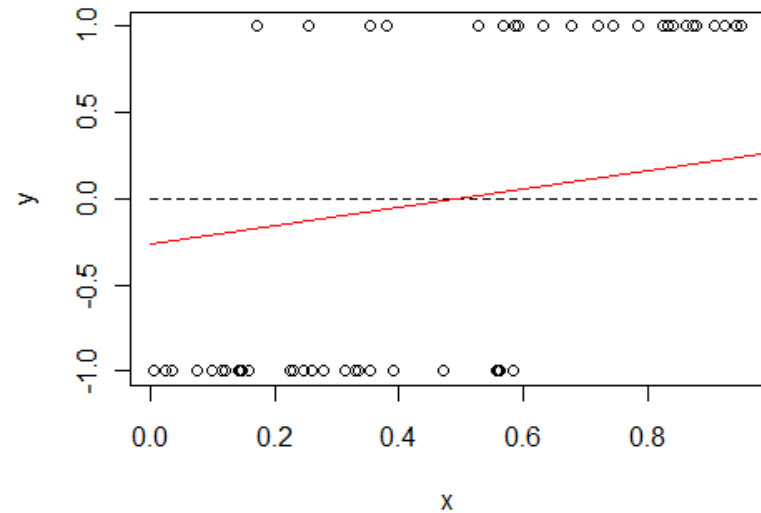
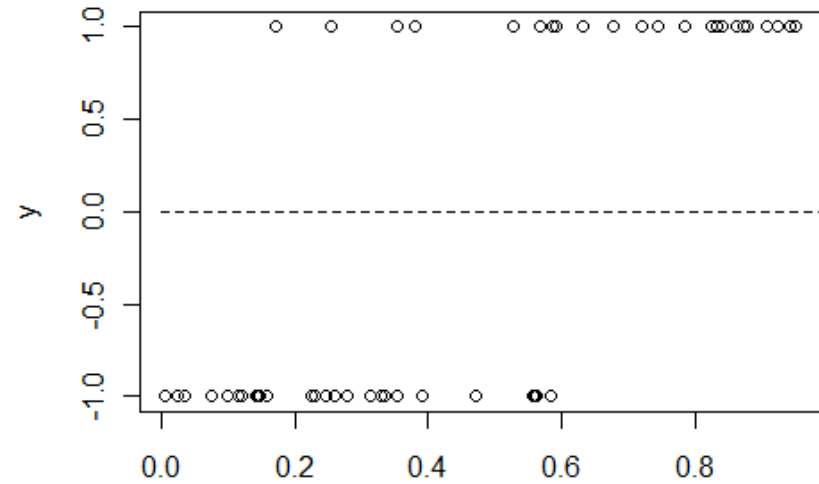
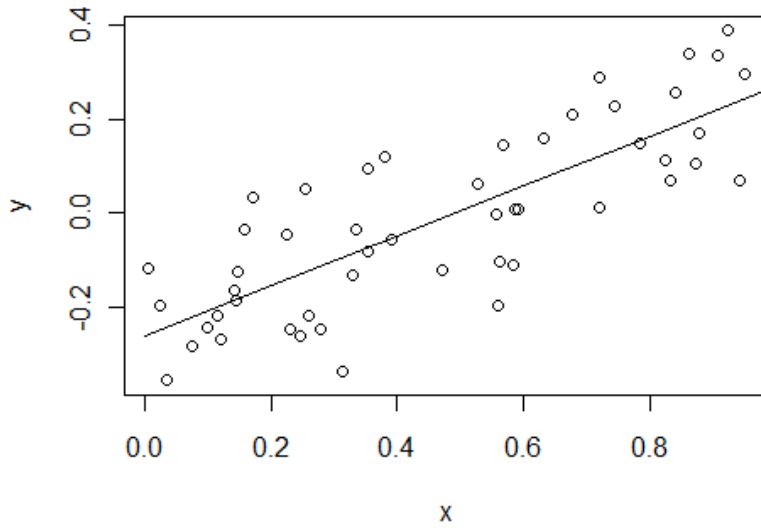
One-dimensional case



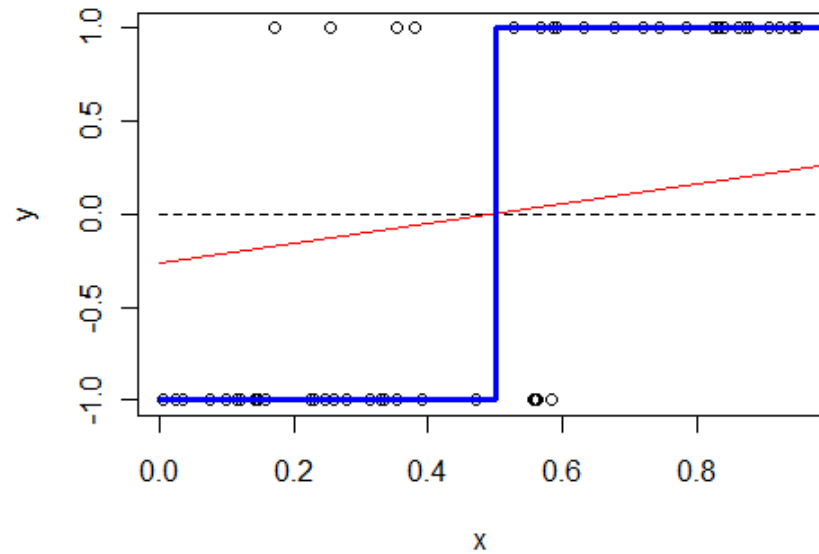
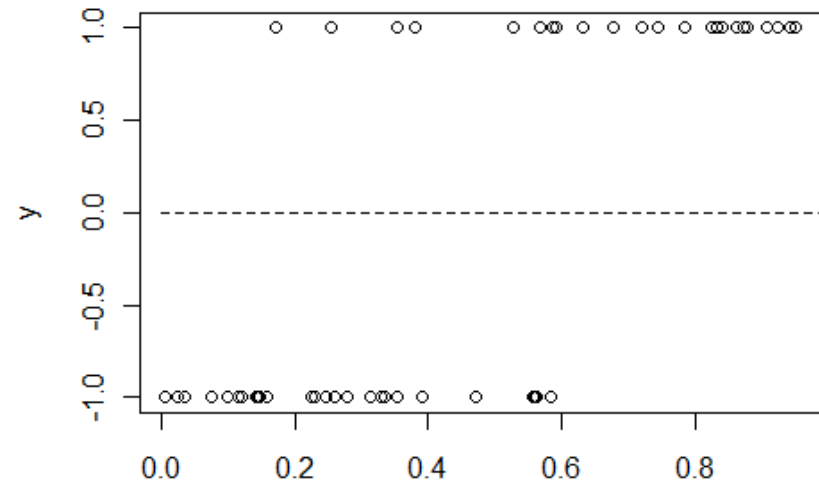
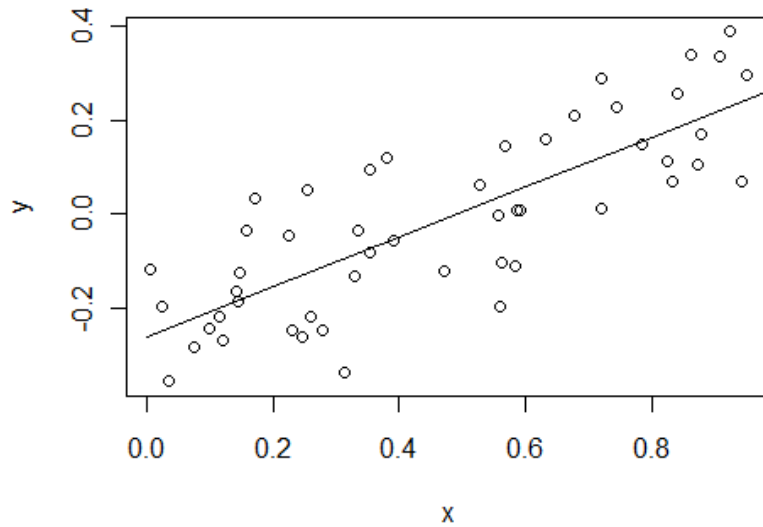
One-dimensional case



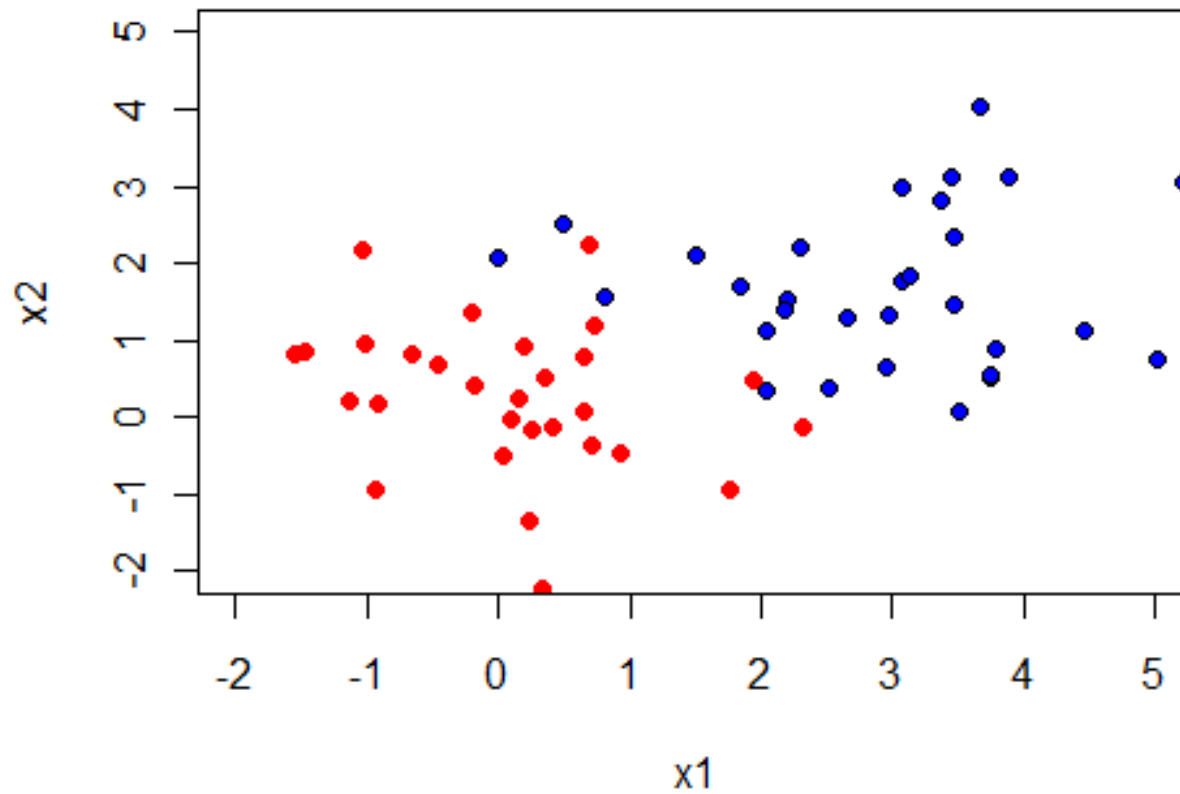
One-dimensional case



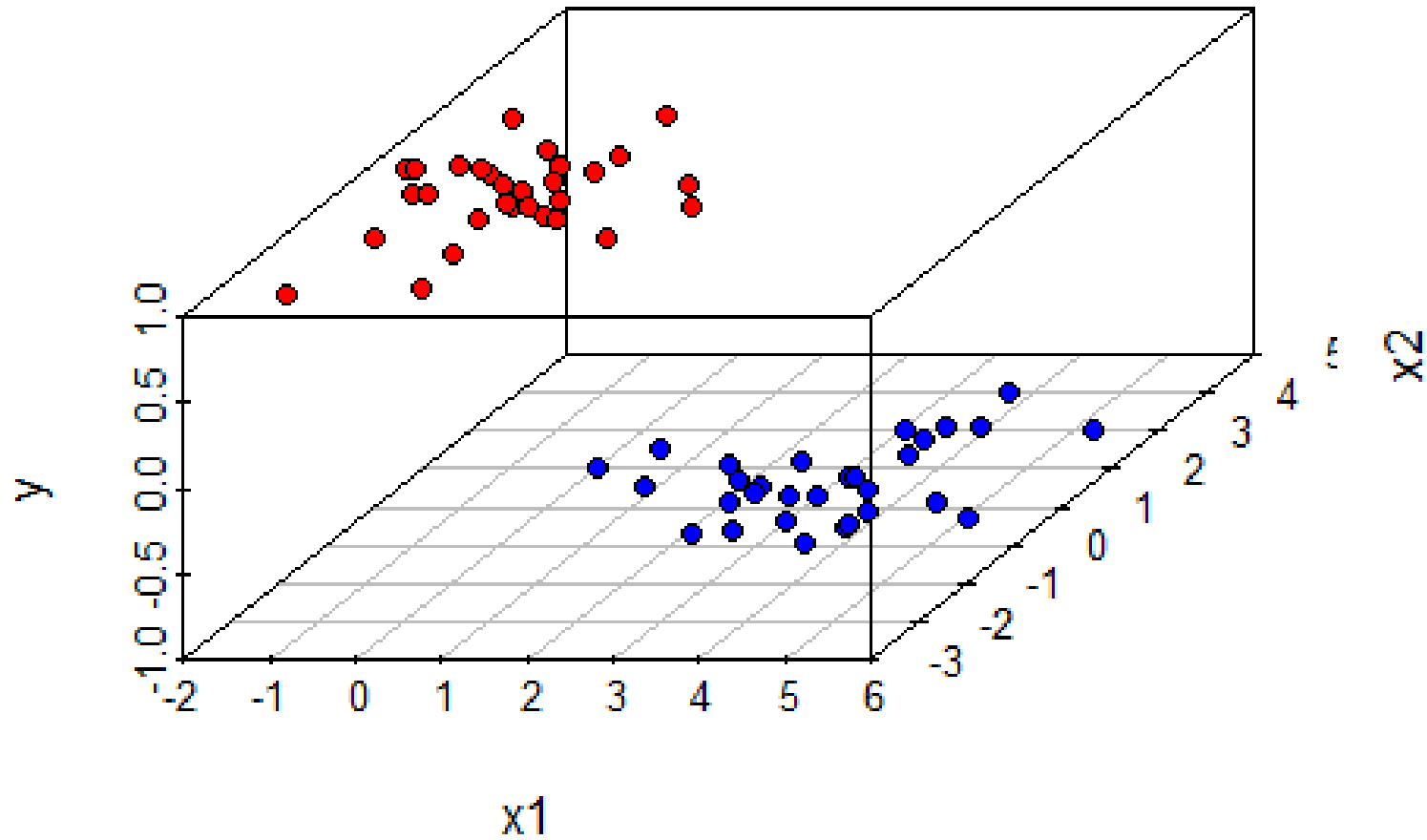
One-dimensional case



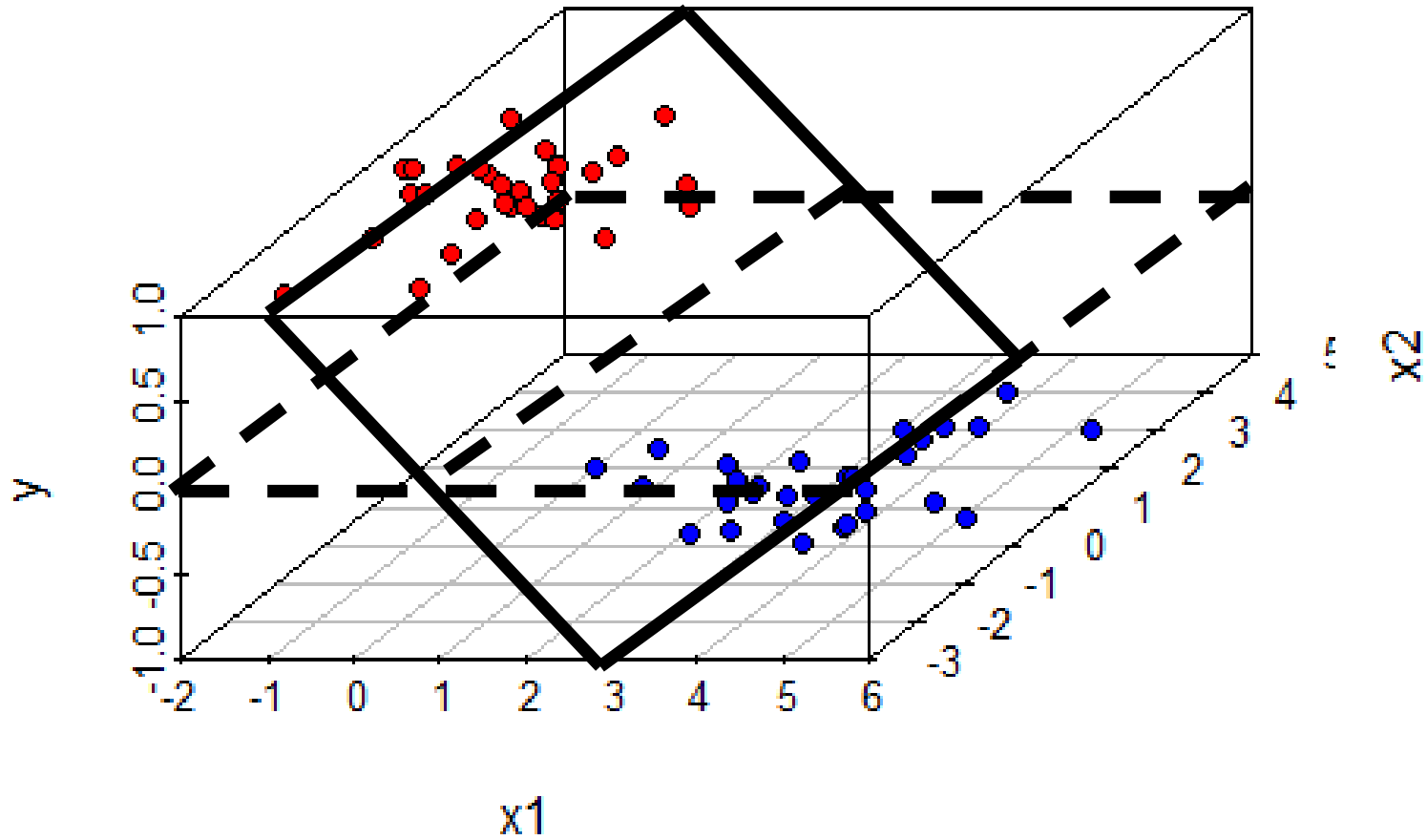
Two-dimensional case



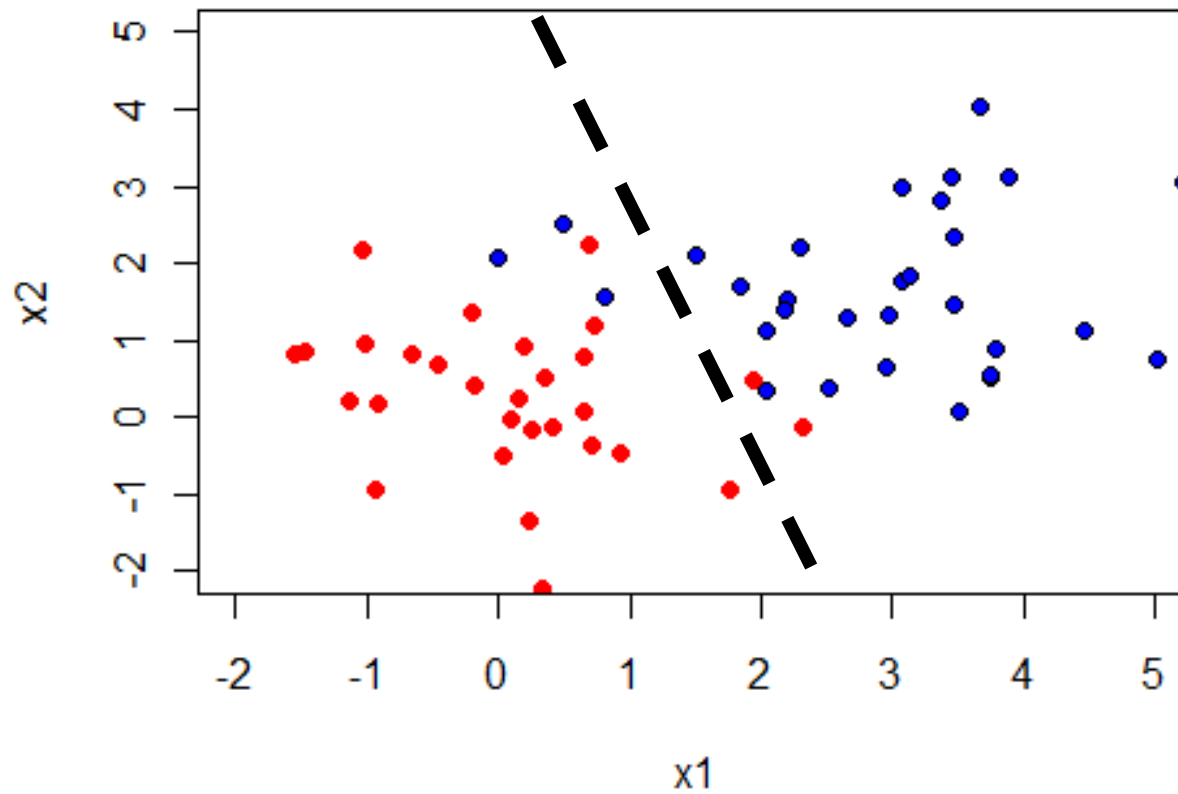
Two-dimensional case



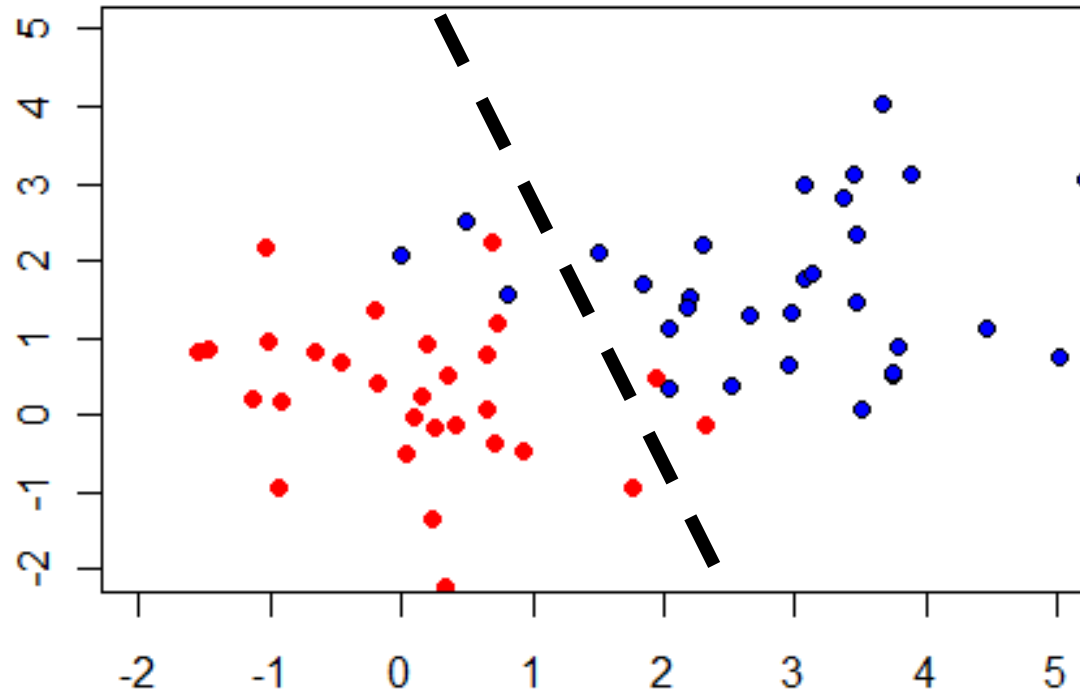
Two-dimensional case



Two-dimensional case

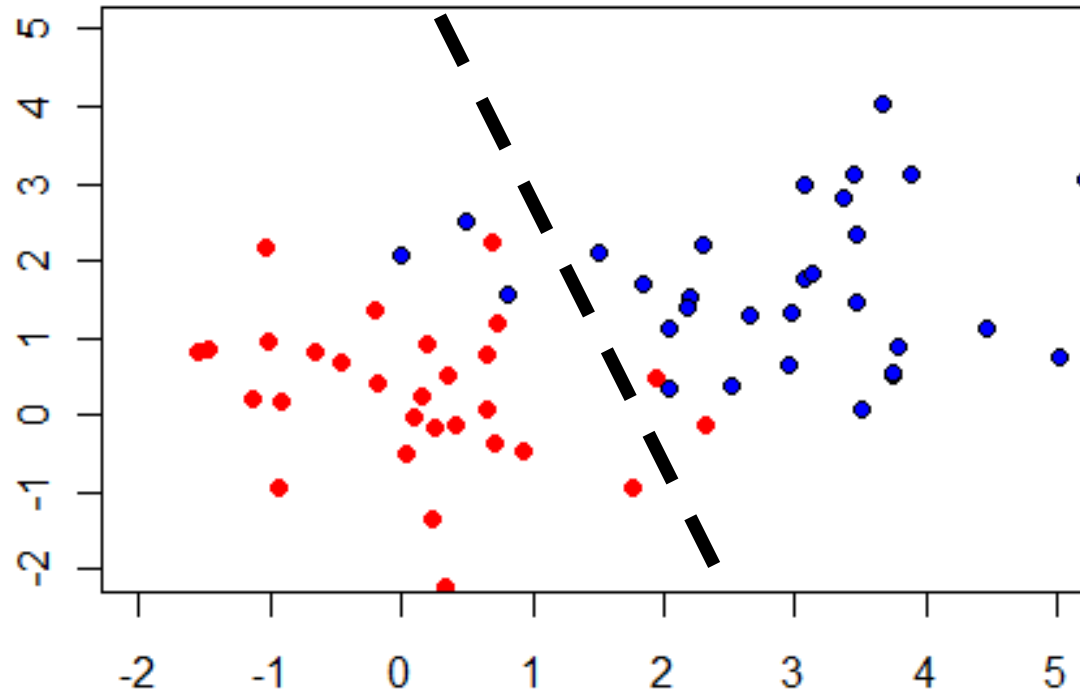


Two-C



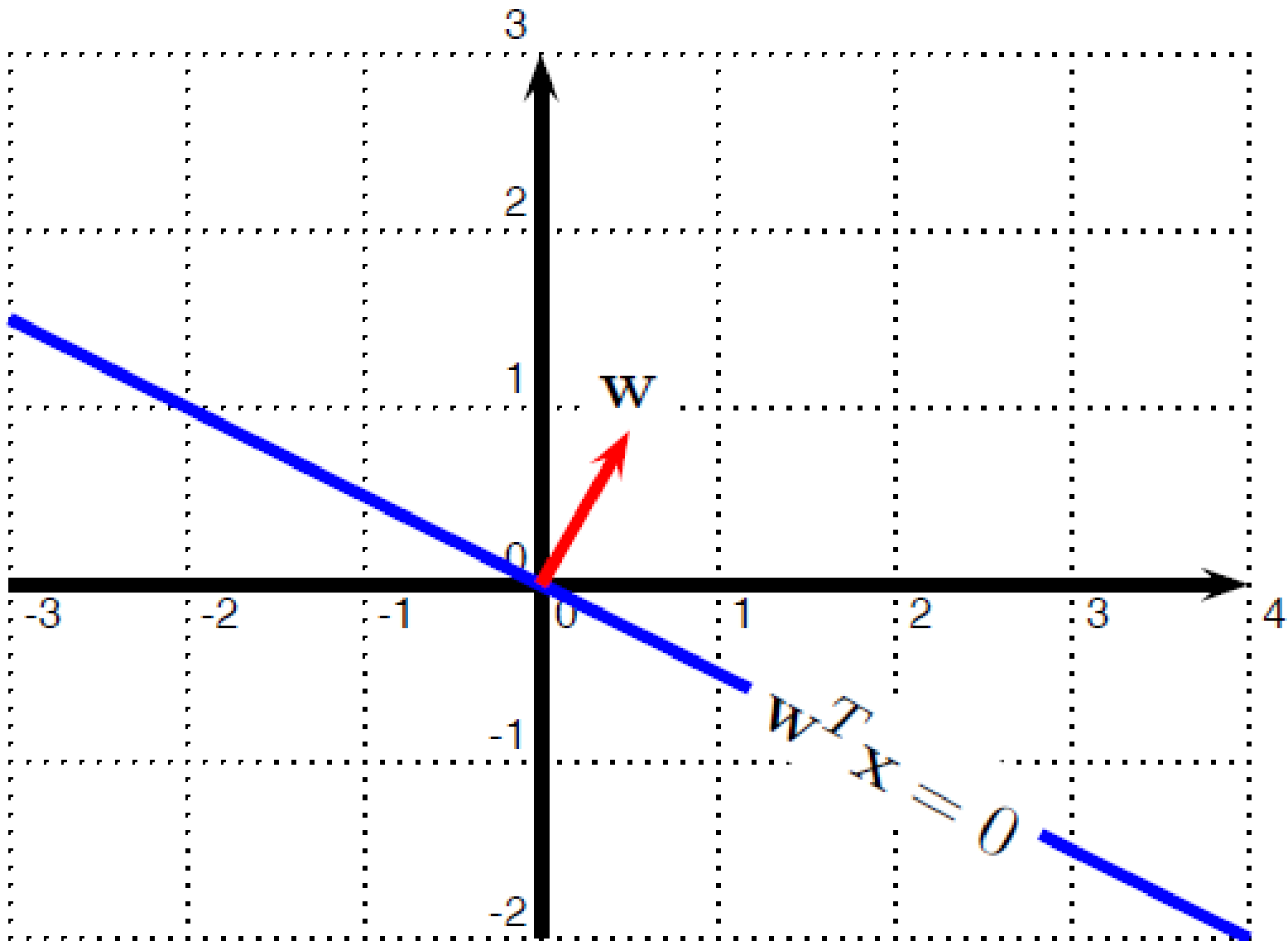
$$w_0 + w_1x_1 + w_2x_2 = 0$$

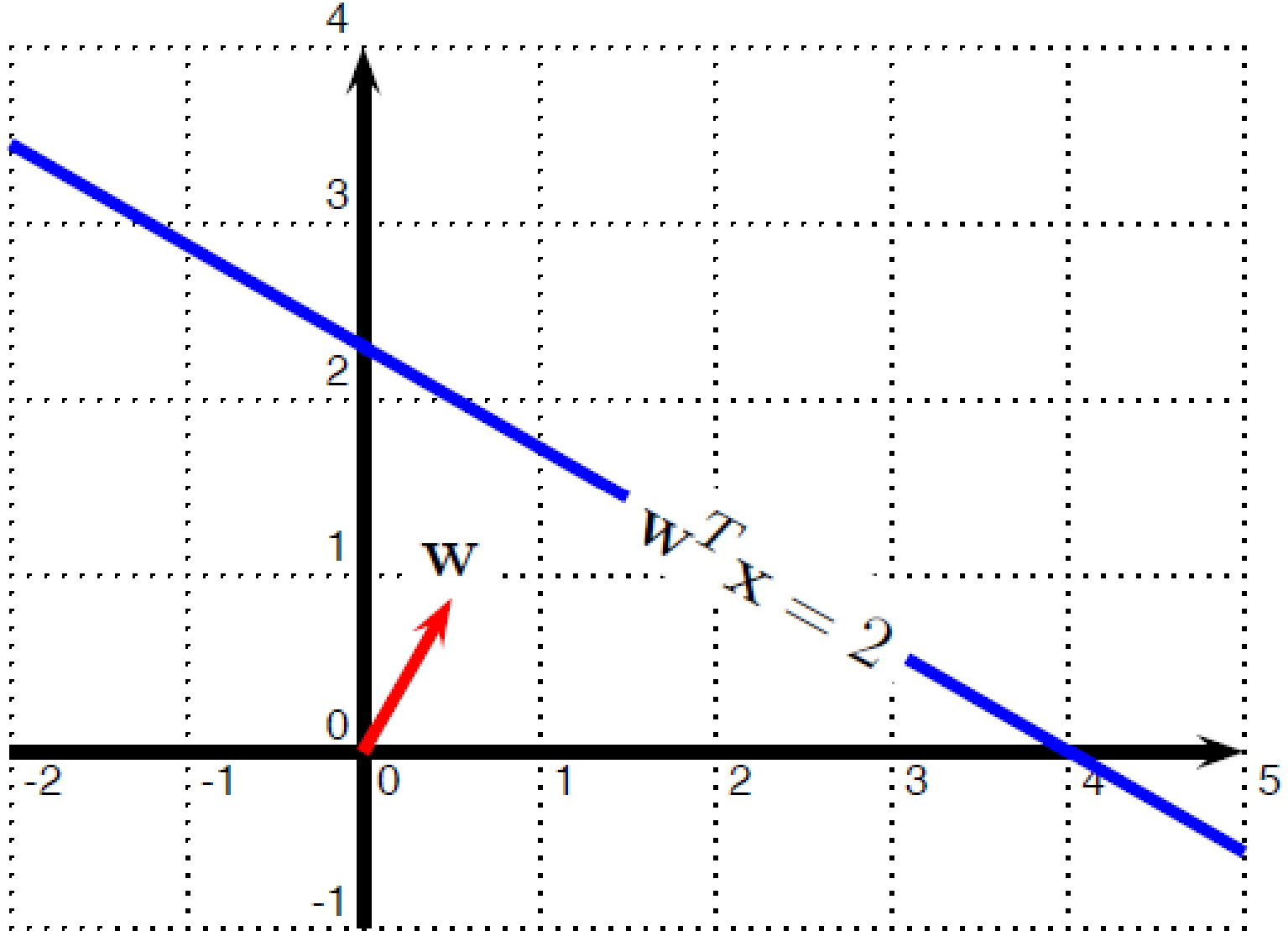
Two-C

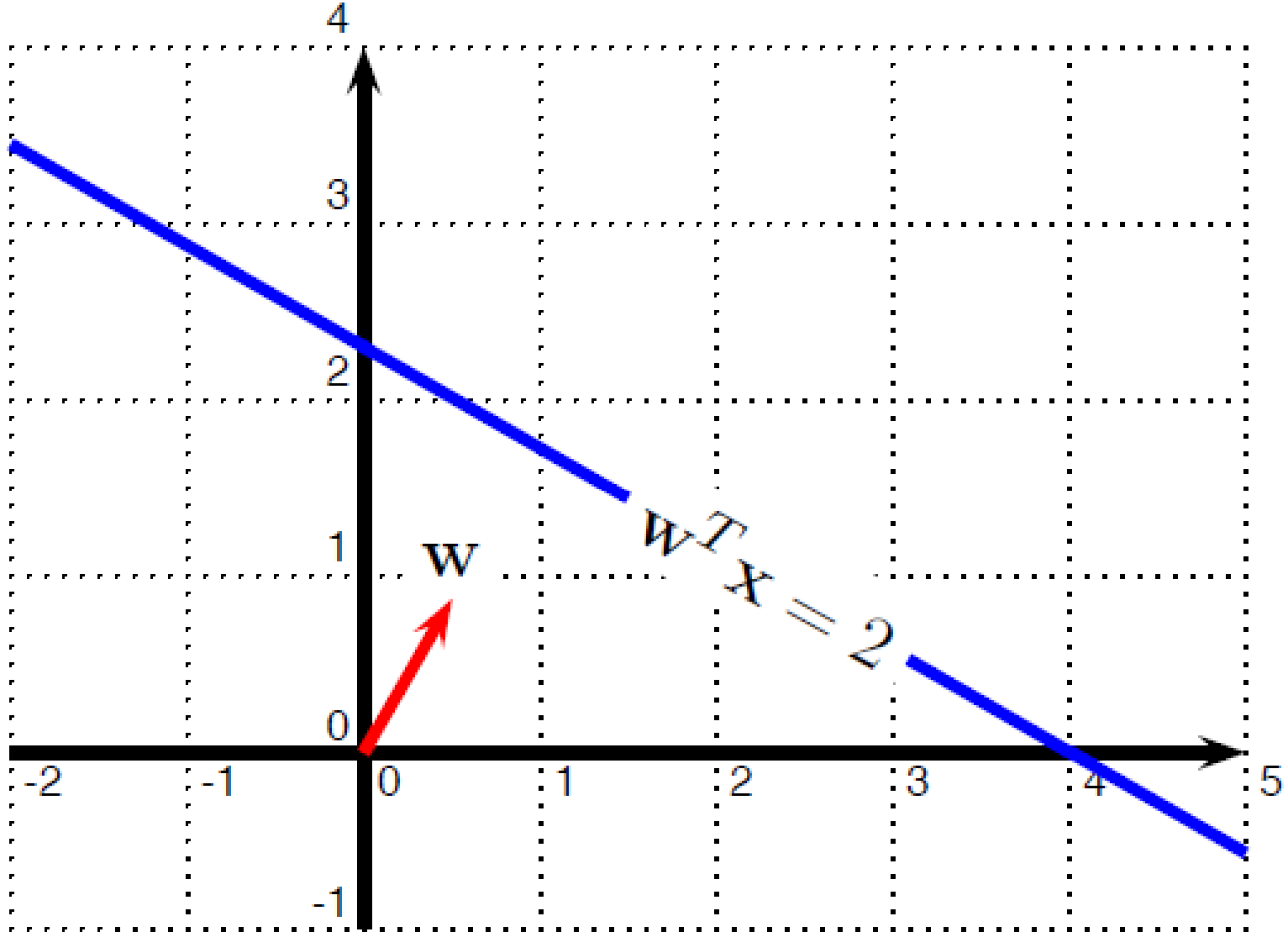


$$w^T x + w_0 = 0$$



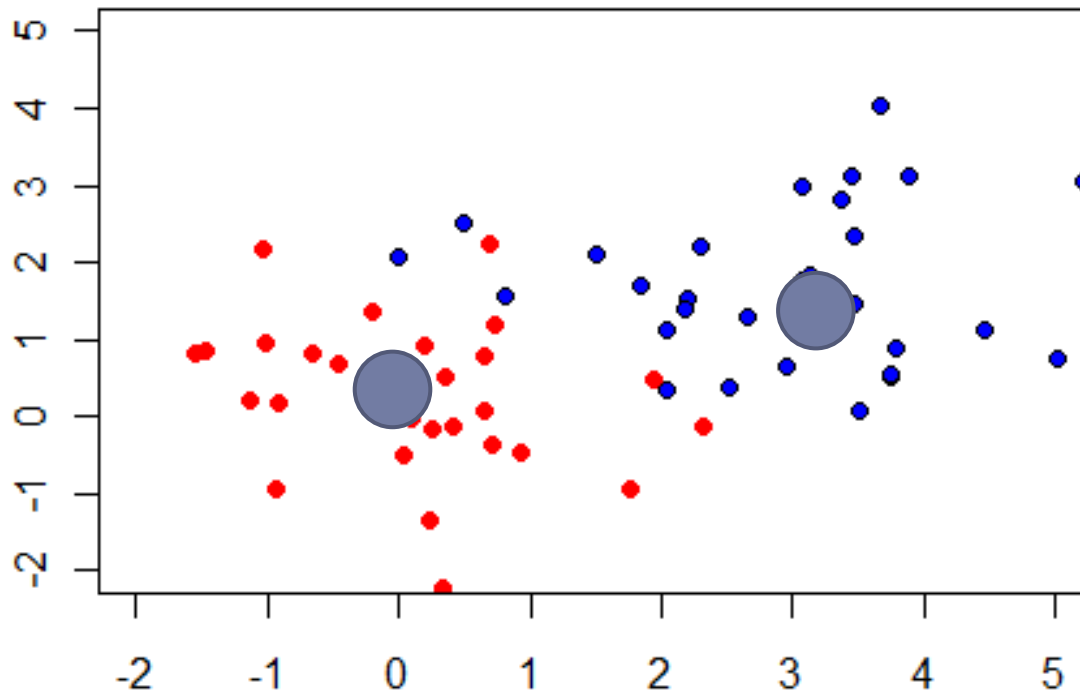




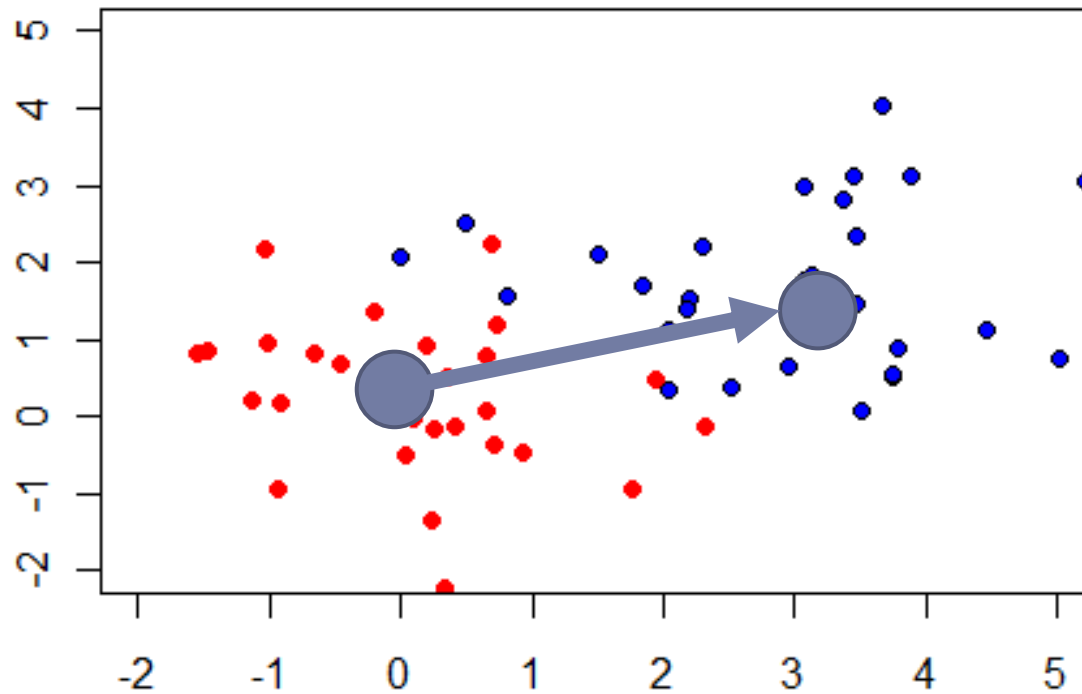


$$w^T (x - p) = 0$$

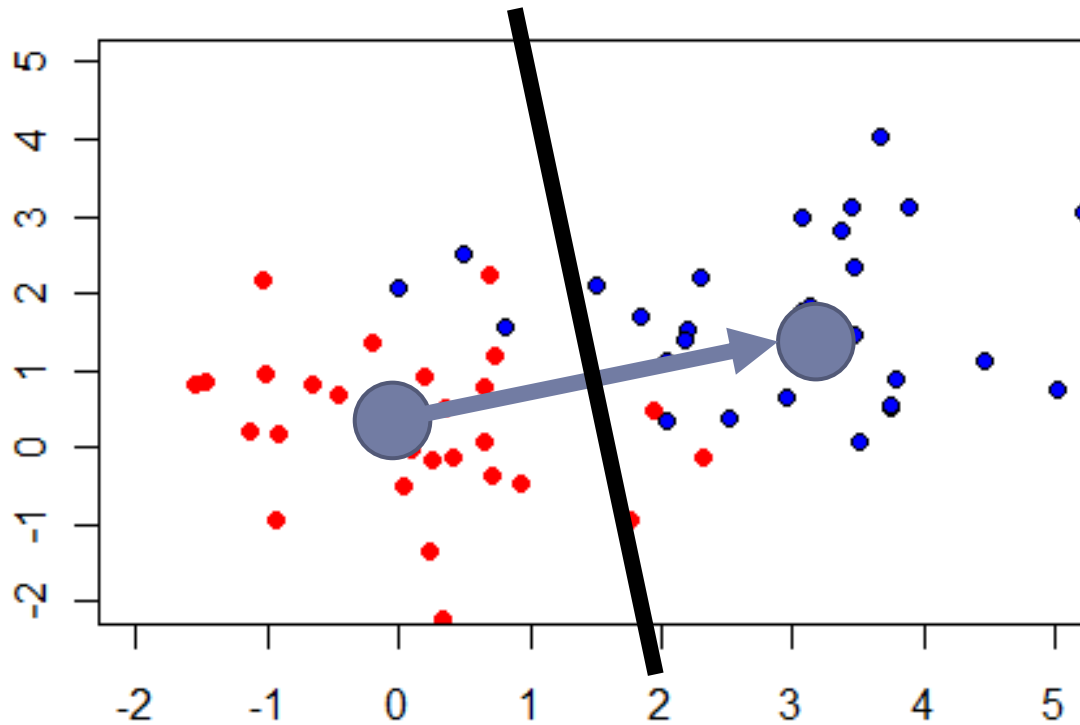
Fisher's discriminant



Fisher's discriminant



Fisher's discriminant



Fisher's discriminant

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

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2.$$

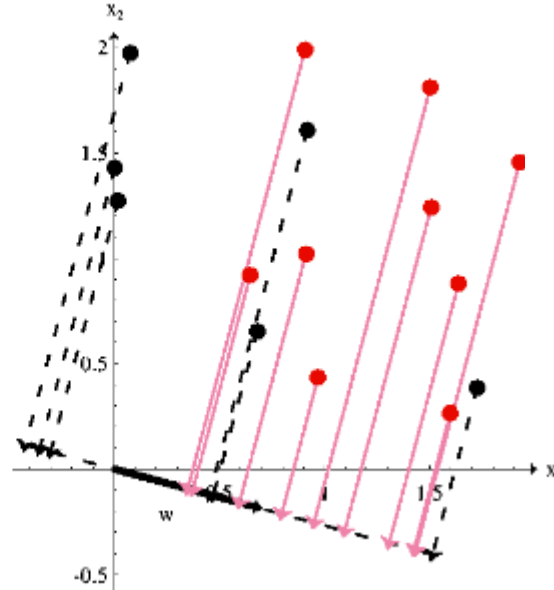
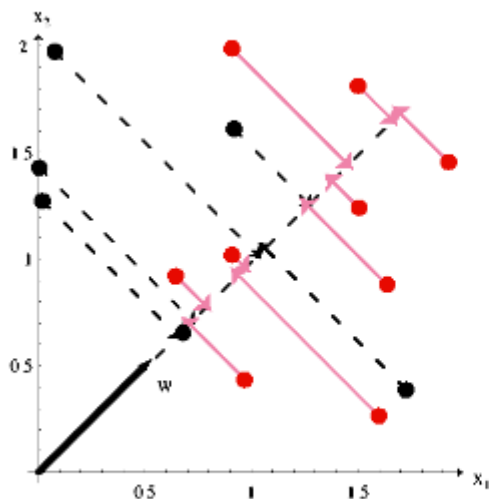
$$\mathbf{S}_i = \sum_{\mathbf{x} \in \mathcal{D}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t$$

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



Interpretations of the FD

- ▶ Probabilistic interpretation
- ▶ 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

- ▶ Recollect linear regression:

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

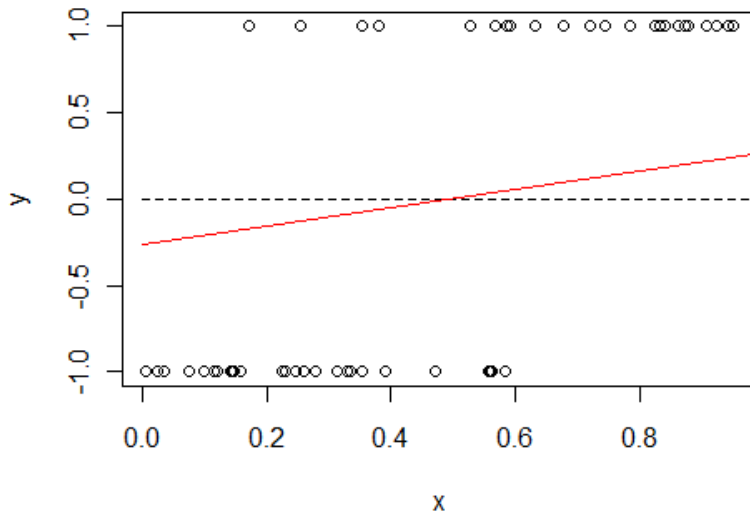
Why not apply the same idea directly to solve the classification tasks?

LSE Methods

- ▶ Recollect linear regression:

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

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



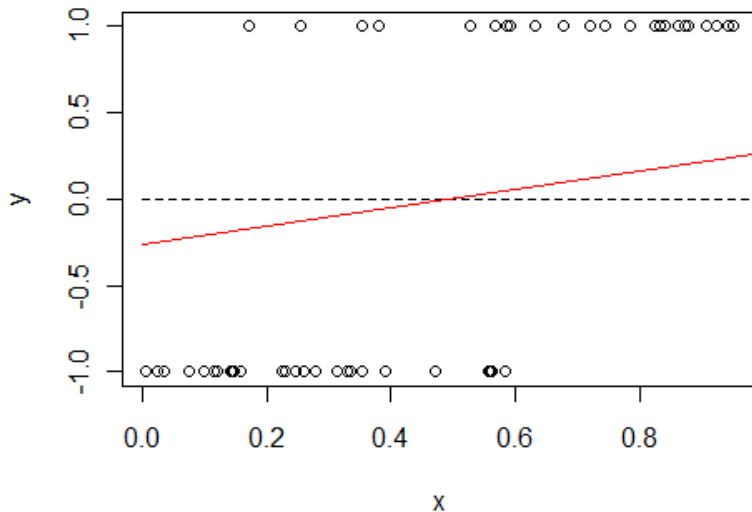
It makes no sense!

LSE Methods

- ▶ Recollect linear regression:

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

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



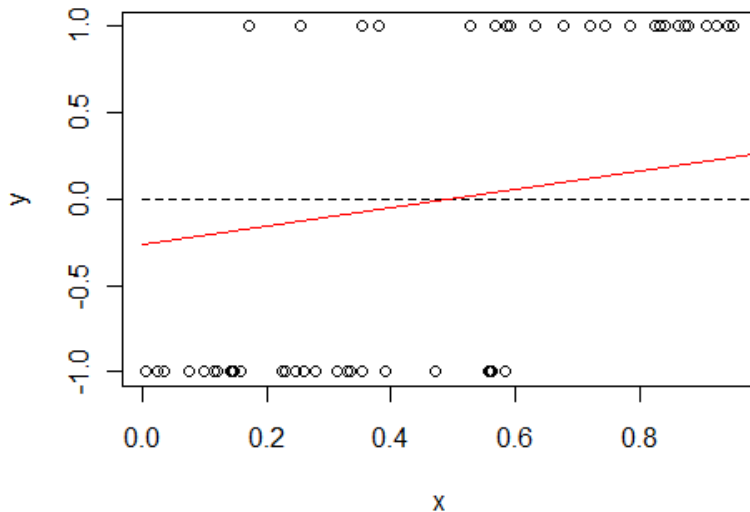
It is too sensitive to outliers

LSE Methods

- ▶ Recollect 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?

- ▶ Obviously, if there exists some good classifier

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

then for some y'_i

such that $y'_i > 0 \iff y_i = 1$, it holds that

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

Hence indeed:

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

LSE Methods – Why?

- ▶ Obviously, if there exists some good classifier

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

then for some y'_i

such that $y'_i > 0 \iff y_i = 1$, it holds that

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

Hence indeed:

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

Unfortunately it is not easy to find the proper \mathbf{y}' .

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.

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

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

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.
- ▶ μ 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?