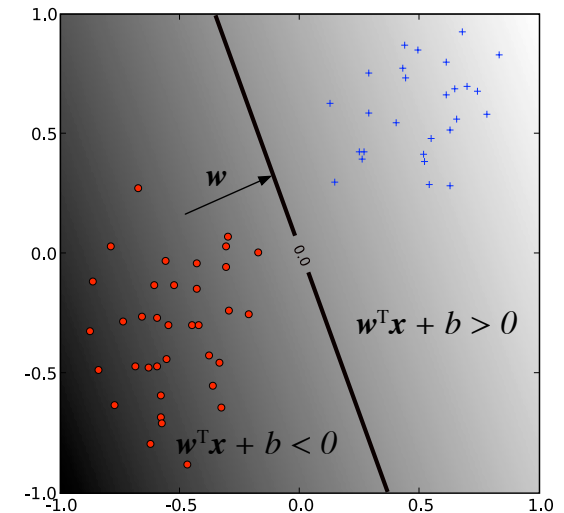

Linear models and the perceptron algorithm

Chapters 1, 3.1



Labeled data

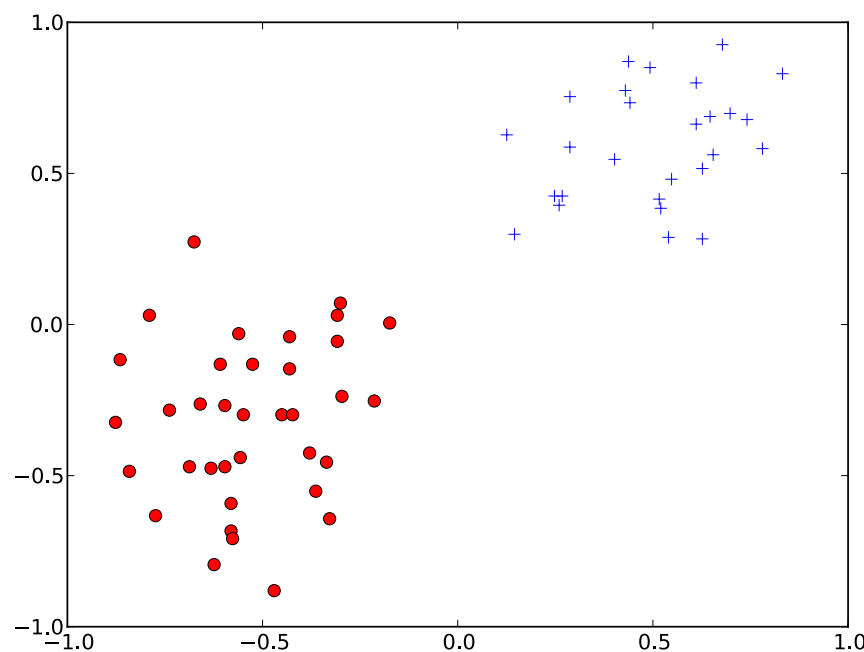
A labeled dataset:

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

Where $\mathbf{x}_i \in \mathbb{R}^d$ are d-dimensional vectors

The labels:

are discrete for classification problems (e.g. +1, -1) for binary classification



Labeled data

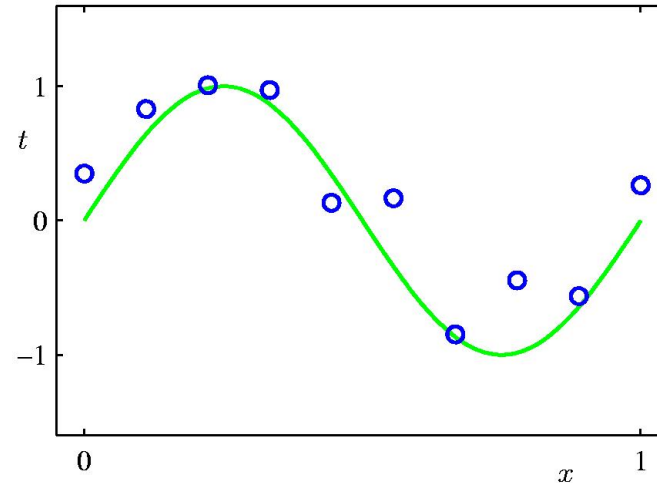
A labeled dataset:

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

Where $\mathbf{x}_i \in \mathbb{R}^d$ are d-dimensional vectors

The labels:

are continuous values for a regression problem



Labeled data

A labeled dataset:

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

Where $\mathbf{x}_i \in \mathbb{R}^d$ are d-dimensional vectors

We typically assume that the examples are i.i.d. (independent and identically distributed), and come from an unknown probability distribution $P(\mathbf{x}, y)$.

Dot products

Definition: The **Euclidean dot product** between two vectors is the expression

$$\mathbf{w}^T \mathbf{x} = \sum_{i=1}^d w_i x_i$$

The dot product is also referred to as inner product or scalar product.

It is sometimes denoted as $\mathbf{w} \cdot \mathbf{x}$
(hence the name dot product).

Dot products

Definition: The **Euclidean dot product** between two vectors is the expression

$$\mathbf{w}^T \mathbf{x} = \sum_{i=1}^d w_i x_i$$

The dot product is also referred to as inner product or scalar product.

Geometric interpretation. The dot product between two unit vectors¹ is the cosine of the angle between them.

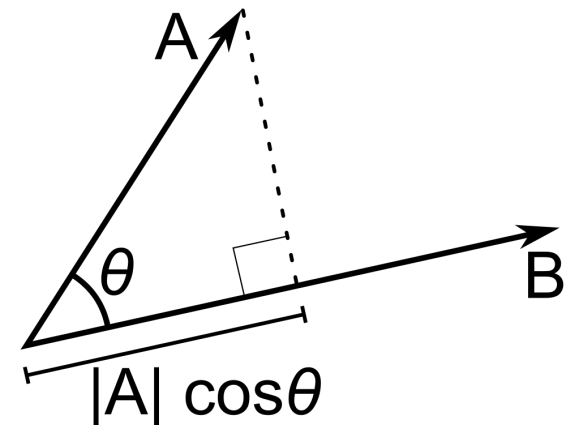
The dot product between a vector and a unit vector is the length of its projection in that direction.

And in general:

$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \cdot \|\mathbf{x}\| \cos(\theta)$$

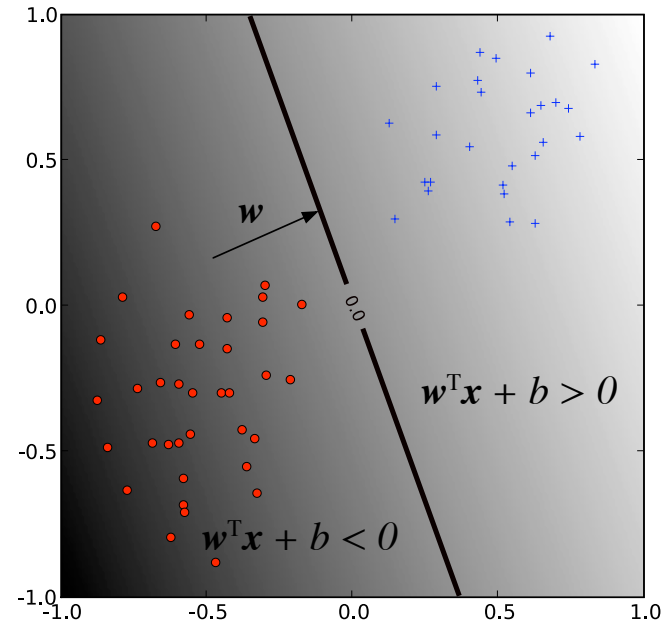
The norm of a vector:

$$\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x}$$

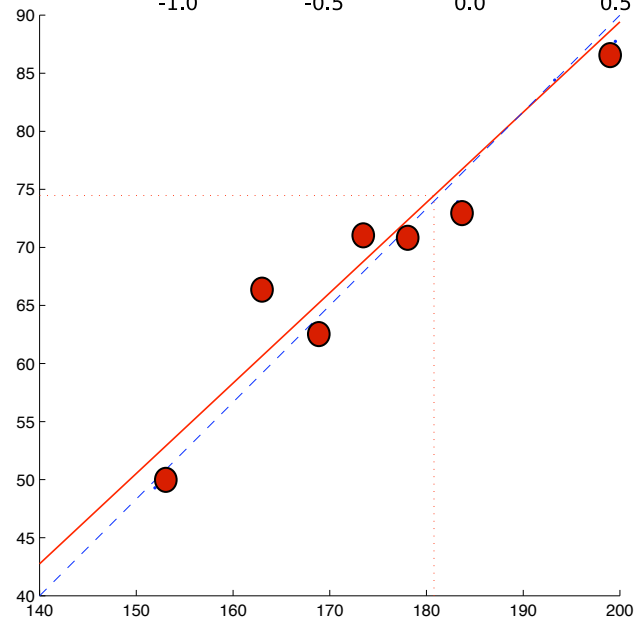


Linear models

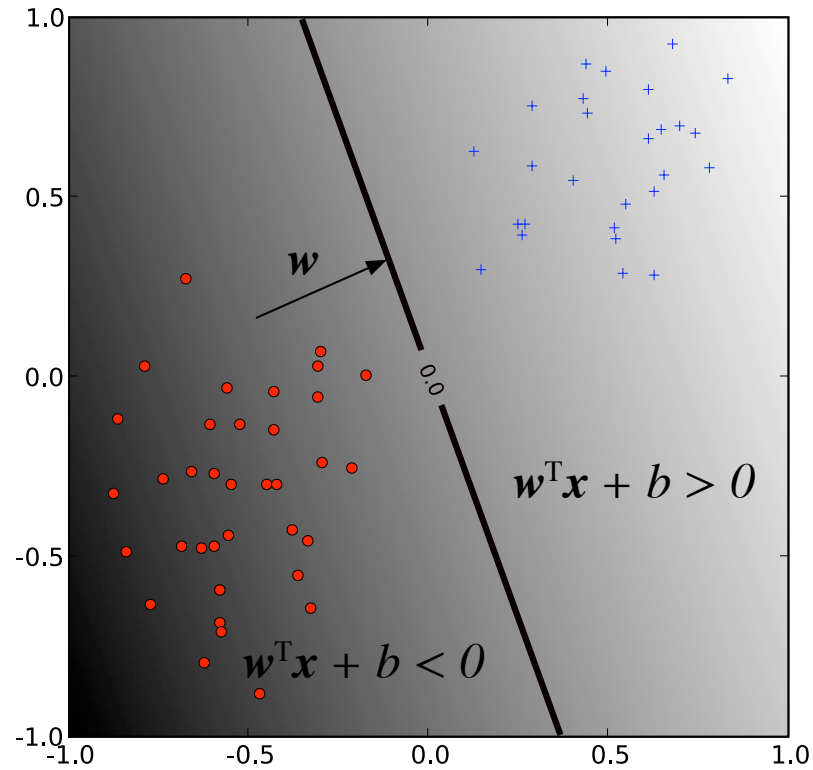
Linear models for classification
(linear decision boundaries)



Linear models for regression
(estimating a linear function)



Linear models for classification

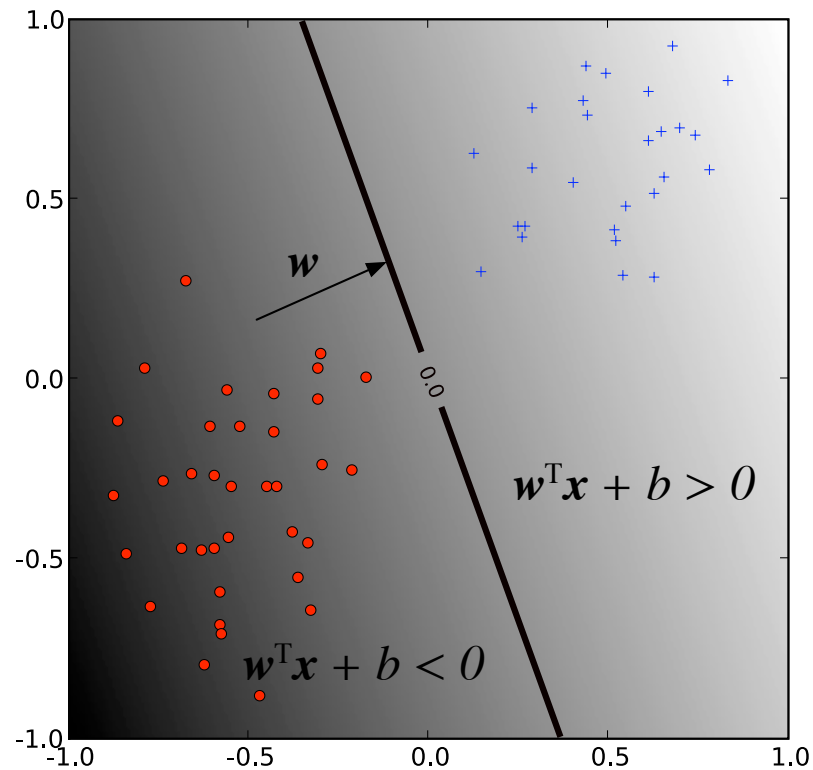


Discriminant/scoring function:

$$w^T x + b$$

weight vector bias

Linear models for classification

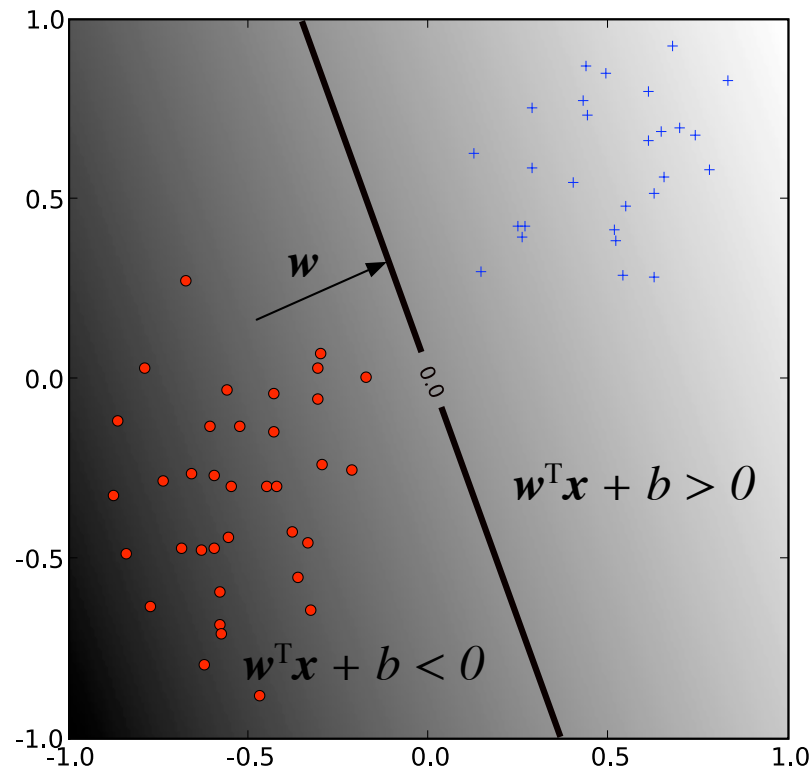


Decision boundary:

all x such that
$$w^T x + b = 0$$

For linear models the the decision boundary is a line in 2-d, a plane in 3-d and a hyperplane in higher dimensions

Linear models for classification

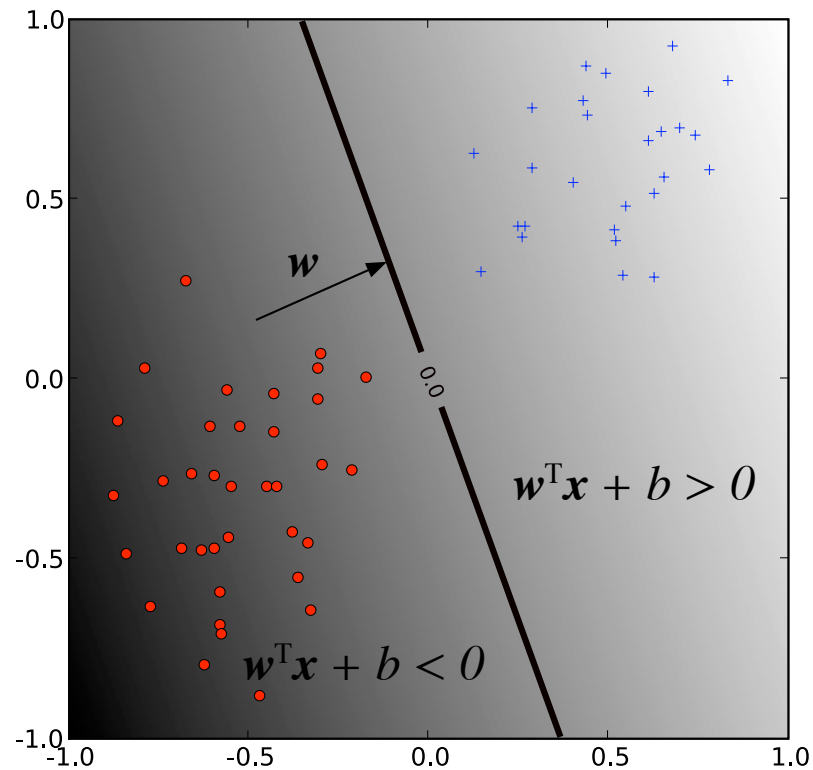


Using the discriminant to
make a prediction:

$$\hat{y} = \text{sign}(w^T x + b)$$

the sign function equals 1 when its argument is positive and -1 otherwise

Linear models for classification



Decision boundary:
all x such that

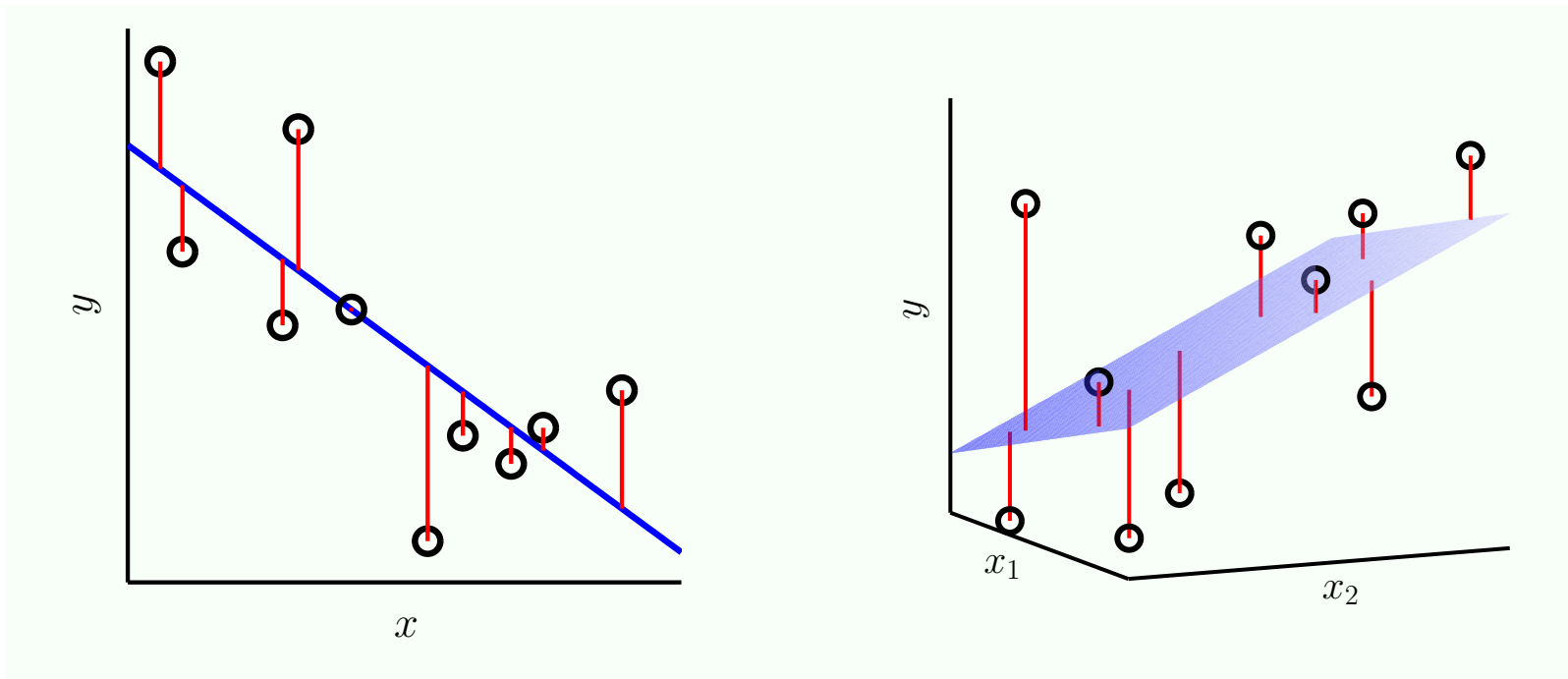
$$w^T x + b = 0$$

What can you say about the decision boundary when $b = 0$?

Linear models for regression

When using a linear model for regression the scoring function is the prediction:

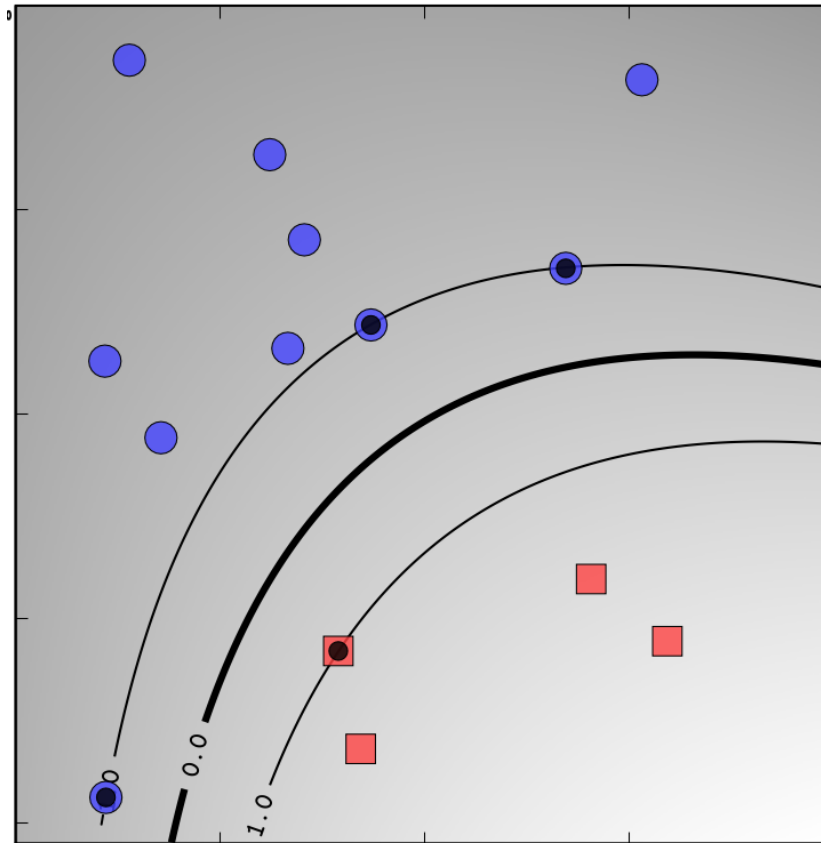
$$\hat{y} = \mathbf{w}^T \mathbf{x} + b$$



Why linear?

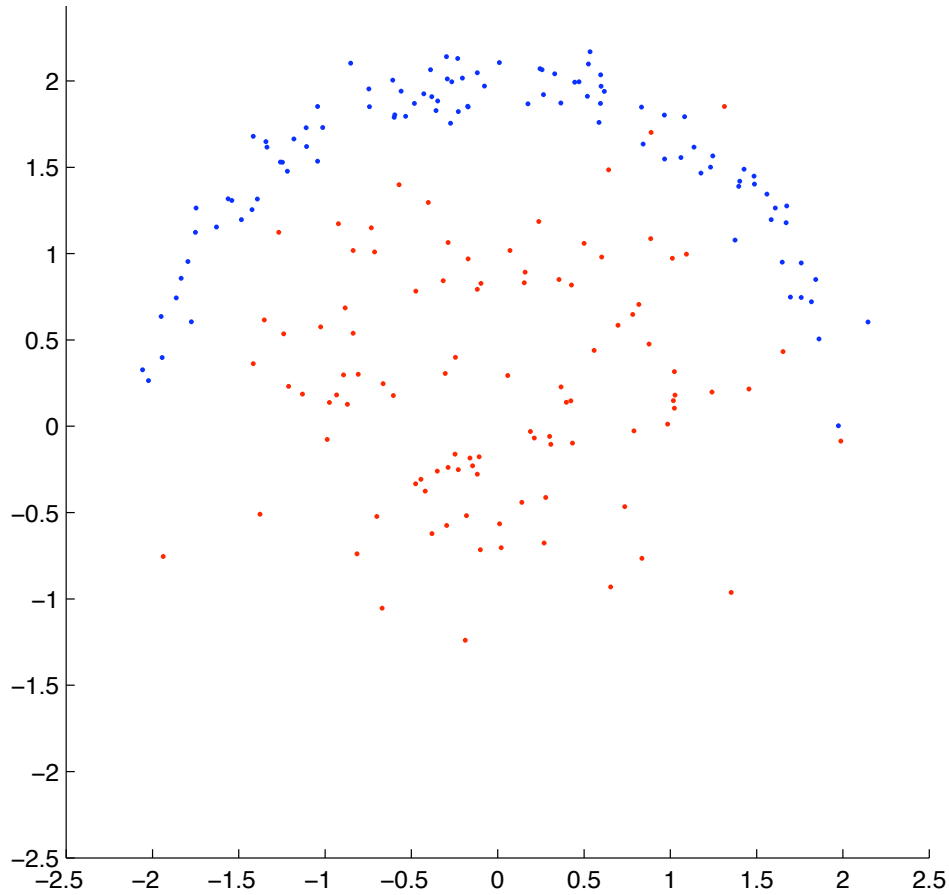
- It's a good baseline: always start simple
- Linear models are stable
- Linear models are less likely to overfit the training data because they have less parameters. Can sometimes underfit. Often all you need when the data is high dimensional.
- Lots of scalable algorithms

From linear to non-linear

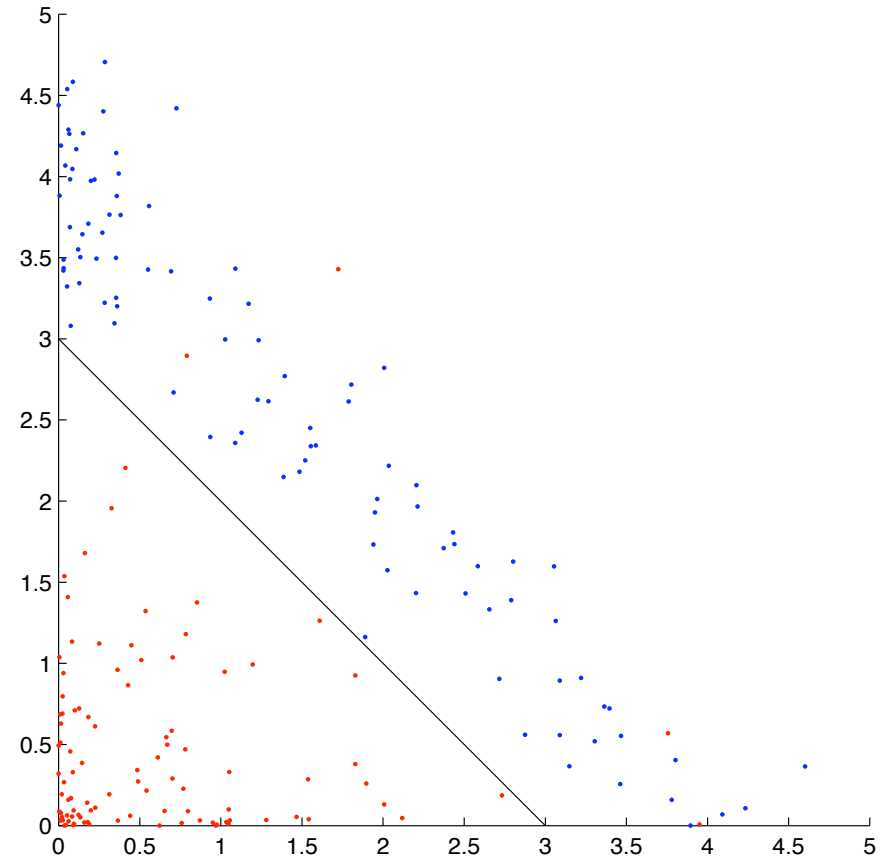


There is a neat mathematical trick that will enable us to use linear classifiers to create non-linear decision boundaries!

From linear to non-linear



Original data: not linearly separable

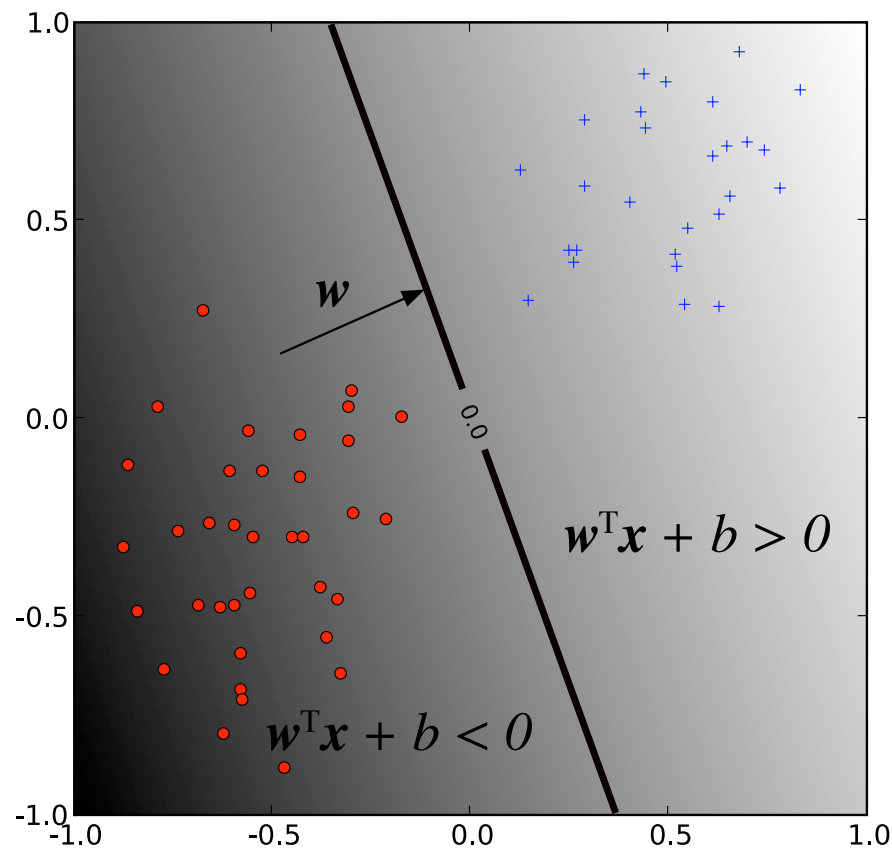


Transformed data: $(x', y') = (x^2, y^2)$

Linearly separable data

Linearly separable data: there exists a linear decision boundary separating the classes.

Example:



The bias and homogeneous coordinates

Formulating a model that does not have a bias does not reduce the expressivity of the model because we can obtain a bias using the following trick:

Add another dimension x_0 to each input and set it to 1.

Learn a weight vector of dimension $d+1$ in this extended space, and interpret w_0 as the bias term. With the notation

$$\mathbf{w} = (w_1, \dots, w_d) \quad \tilde{\mathbf{w}} = (w_0, w_1, \dots, w_d)$$

$$\tilde{\mathbf{x}} = (1, x_1, \dots, x_d)$$

We have that:
$$\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}} = w_0 + \mathbf{w}^\top \mathbf{x}$$

Finding a good hyperplane

We would like a classifier that fits the data, i.e. we would like to find a vector \mathbf{w} that minimizes

$$\begin{aligned} E_{\text{in}} &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}[h(\mathbf{x}_i) \neq f(\mathbf{x}_i)] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}[\text{sign}(\mathbf{w}^\top \mathbf{x}_i) \neq y_i] \end{aligned}$$

This is a difficult problem because of the discrete nature of the indicator and sign function (known to be NP-hard).

The perceptron algorithm (Rosenblatt, 1957)

Idea: iterate over the training examples, and update the weight vector \mathbf{w} in a way that would make \mathbf{x}_i is more likely to be correctly classified.

Let's assume that \mathbf{x}_i is misclassified, and is a positive example i.e.

$$\mathbf{w}^T \mathbf{x}_i < 0$$

Note: we're learning a classifier without a bias term

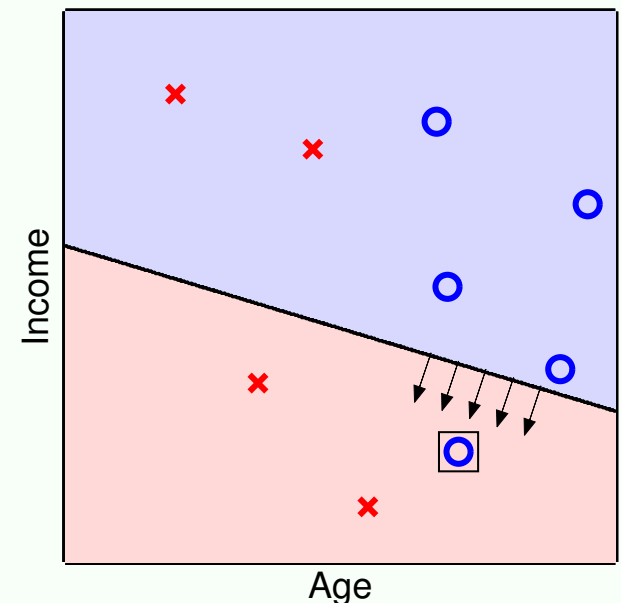
We would like to update \mathbf{w} to \mathbf{w}' such that

$$(\mathbf{w}')^T \mathbf{x}_i > \mathbf{w}^T \mathbf{x}_i$$

This can be achieved by choosing

$$\mathbf{w}' = \mathbf{w} + \eta \mathbf{x}_i$$

$0 < \eta \leq 1$ is the learning rate



Rosenblatt, Frank (1957), The Perceptron--a perceiving and recognizing automaton. Report 85-460-1, Cornell Aeronautical Laboratory.

Section 1.1 in the book

The perceptron algorithm

If \mathbf{x}_i is a negative example, the update needs to be opposite.

Overall, we can summarize the two cases as:

$$\mathbf{w}' = \mathbf{w} + \eta y_i \mathbf{X}_i$$

Input: labeled data D in homogeneous coordinates

Output: weight vector \mathbf{w}

```
 $\mathbf{w} = 0$ 
```

```
converged = false
```

```
while not converged :
```

```
    converged = true
```

```
    for  $i$  in  $1, \dots, N$  :
```

```
        if  $\mathbf{x}_i$  is misclassified update  $\mathbf{w}$  and set
```

```
            converged=false
```

```
return  $\mathbf{w}$ 
```

Demo

Show jupyter notebook on the perceptron

The perceptron algorithm

Since the algorithm is not guaranteed to converge if the data is not linearly separable you need to set a limit on the number of iterations:

Input: labeled data D in homogeneous coordinates

Output: weight vector \mathbf{w}

```
 $\mathbf{w} = 0$ 
```

```
converged = false
```

```
while (not converged or number of iterations < T) :
```

```
    converged = true
```

```
    for i in 1,...,N :
```

```
        if  $\mathbf{x}_i$  is misclassified:
```

```
            update  $\mathbf{w}$  and set converged=false
```

```
return  $\mathbf{w}$ 
```

The perceptron algorithm

The algorithm is guaranteed to converge if the data is linearly separable, and does not converge otherwise.

Issues with the algorithm:

- The algorithm chooses an **arbitrary** hyperplane that separates the two classes. It may not be the best one from the learning perspective.
- Does not converge if the data is not separable (can halt after a fixed number of iterations).

There are variants of the algorithm that address these issues (to some extent).

The pocket algorithm

Input: labeled data D in homogeneous coordinates

Output: a weight vector \mathbf{w}

$\mathbf{w} = \mathbf{0}$, $\mathbf{w}_{\text{pocket}} = \mathbf{0}$

converged = false

while (not converged or number of iterations $< T$) :

 converged = true

 for i in $1, \dots, N$:

 if \mathbf{x}_i is misclassified:

 update \mathbf{w} and set converged=false

 if \mathbf{w} leads to better E_{in} than $\mathbf{w}_{\text{pocket}}$:

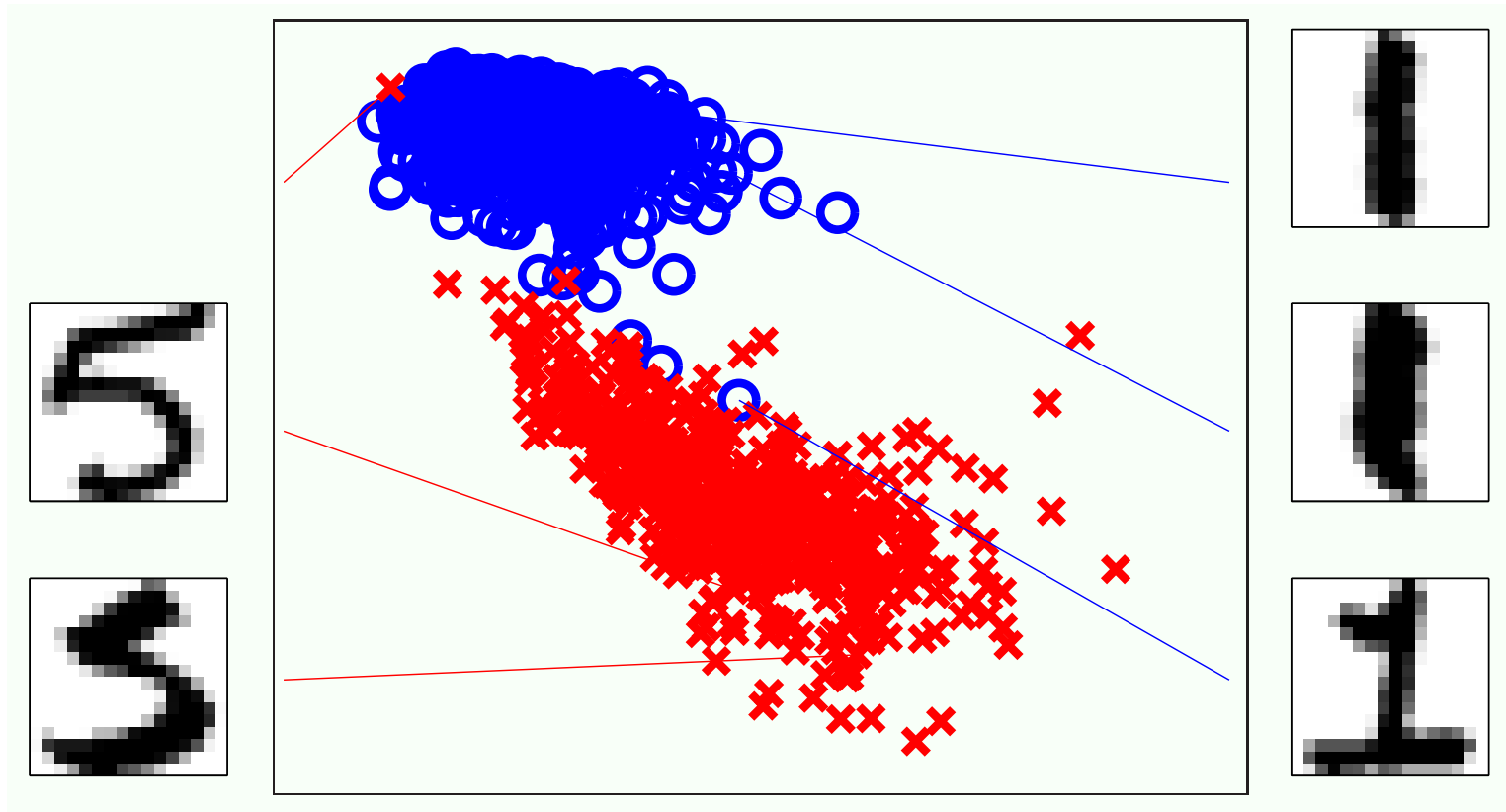
$\mathbf{w}_{\text{pocket}} = \mathbf{w}$

return $\mathbf{w}_{\text{pocket}}$

Gallant, S. I. (1990). Perceptron-based learning algorithms. IEEE Transactions on Neural Networks, vol. 1, no. 2, pp. 179-191.

image classification

Features: important properties of the input you think are relevant for classification



In this case we consider the level of symmetry (image - its flipped version) and overall intensity (fraction of pixels that are dark)

pocket algorithm vs perceptron

Comparison on image data: distinguishing between the digits "1" and "5" (see page 83 in the book):

