

Practical Introduction to Machine Learning

Dăscălescu Dana, March 2023

Scope

Objectives

- Guide participants through the steps required to solve real-world problems using machine learning
- Provide hands-on experience working with real-world data sets and popular ML tools

Focus

- Not about defining the field of ML or providing a deep understanding of the mathematical foundation of ML algorithms
- Focuses on guiding participants through the steps required to solve real-world problems using machine learning

Overview

Data preprocessing

Bag-of-Words:

- *A method of representing text data based on the frequency of occurrence words within documents*

Standardization

Model selection and evaluation

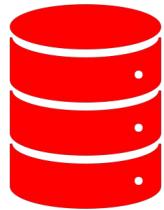
Support Vector Machines

Optimization

Finding the best performing model:

- Grid search

Problem definition



Data

MOROCO –The Moldavian and Romanian Dialectal Corpus
Text collected from the news domain

The samples belong to one of the following six topics: (0) culture, (1) finance, (2) politics, (3) science, (4) sports, (5) tech



Process

Prepare the data: transform the text data into a numerical representation

Pre-process the data: e.g. normalization

Train a classifier

Optimize hyperparameters



Output

For an unseen text document, predict the main category to which it belongs

Data preprocessing – Bag-of-words model

- This is the first sentence
- This sentence is the second sentence
- And this is the third sentence
- Is this the first sentence

Data preprocessing – Bag-of-words model

- This is the first sentence this is the first sentence second and third
- This sentence is the second sentence
- And this is the third sentence
- Is this the first sentence

Data preprocessing – Bag-of-words model

- This is the first sentence
- This sentence is the second sentence
- And this is the third sentence
- Is this the first sentence

this	is	the	first	sentence	second	and	third
1	1	1	1	1	0	0	0
1	1	1	0	2	1	0	0
1	1	1	0	1	0	1	1
1	1	1	1	1	0	0	0

Data preprocessing – Bag-of-words model

- This is the first sentence
- This sentence is the second sentence
- And this is the third sentence
- Is this the first sentence

this is the first sentence second and third

Corpus



feature

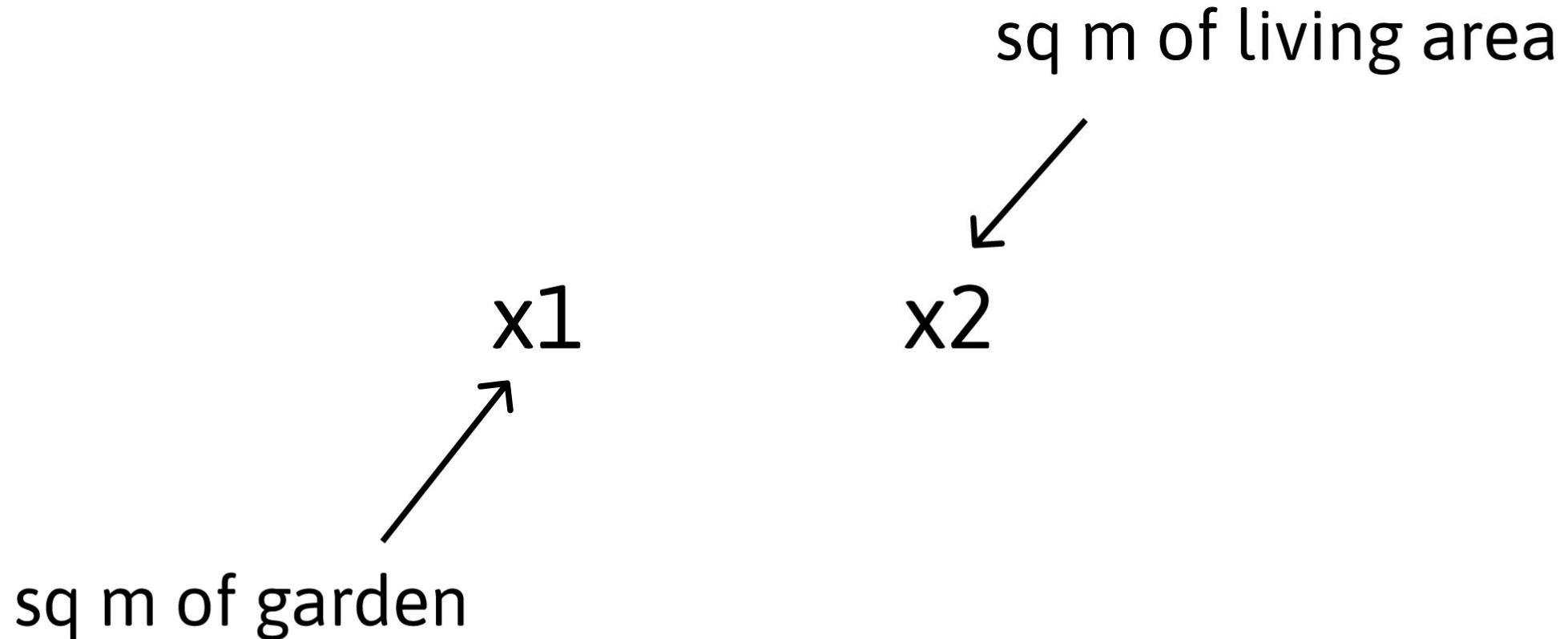
feature 5 of sample 2

1	1	1	1	1	0	0	0
1	1	1	0	2	1	0	0
1	1	0	1	0	1	1	1
1	1	1	1	1	0	0	0
1	1	1	1	1	0	0	0

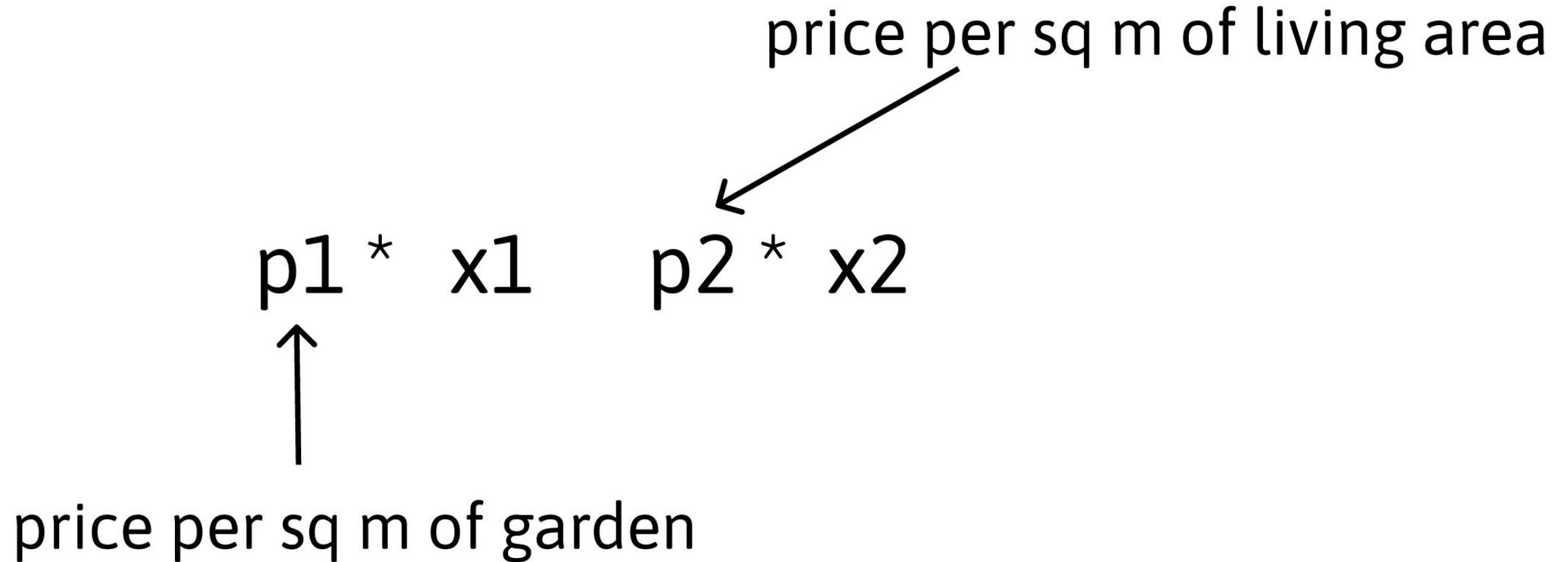
How do we perform classification?

Let's take a more practical example

Price house prediction



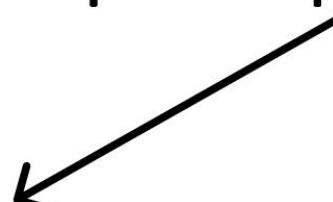
Price house prediction



Price house prediction

$$p_1 * x_1 + p_2 * x_2$$

price per sq m of living area



Price house prediction

Decide whether is an expensive house or not?

$$p_1 * x_1 + p_2 * x_2 \stackrel{>}{\leq} t$$



threshold

Price house prediction

Decide whether is an expensive house or not?

$$p_1 * x_1 + p_2 * x_2 - t \stackrel{>}{\leqslant} 0$$

Price house prediction

Decide whether is an expensive house or not?

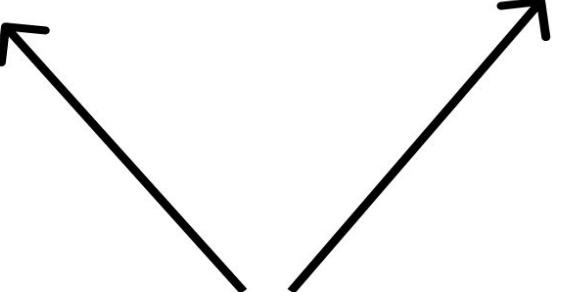
$$p_1 * x_1 + p_2 * x_2 + t \stackrel{>}{\leq} 0$$



sign doesn't matter

Price house prediction

Decide whether is an expensive house or not?

$$w_1 * x_1 + w_2 * x_2 + t \stackrel{>}{\leq} 0$$


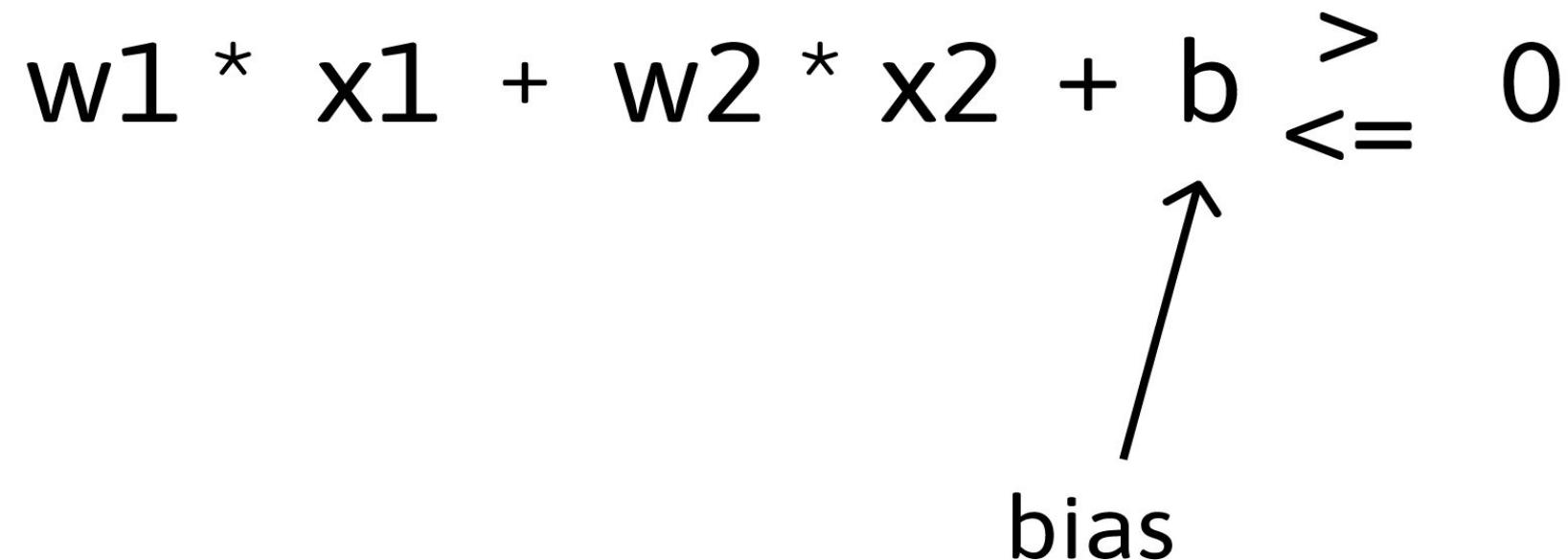
weights

Price house prediction

Decide whether is an expensive house or not?

$$w_1 * x_1 + w_2 * x_2 + b \begin{matrix} > \\ \leq \end{matrix} 0$$

bias

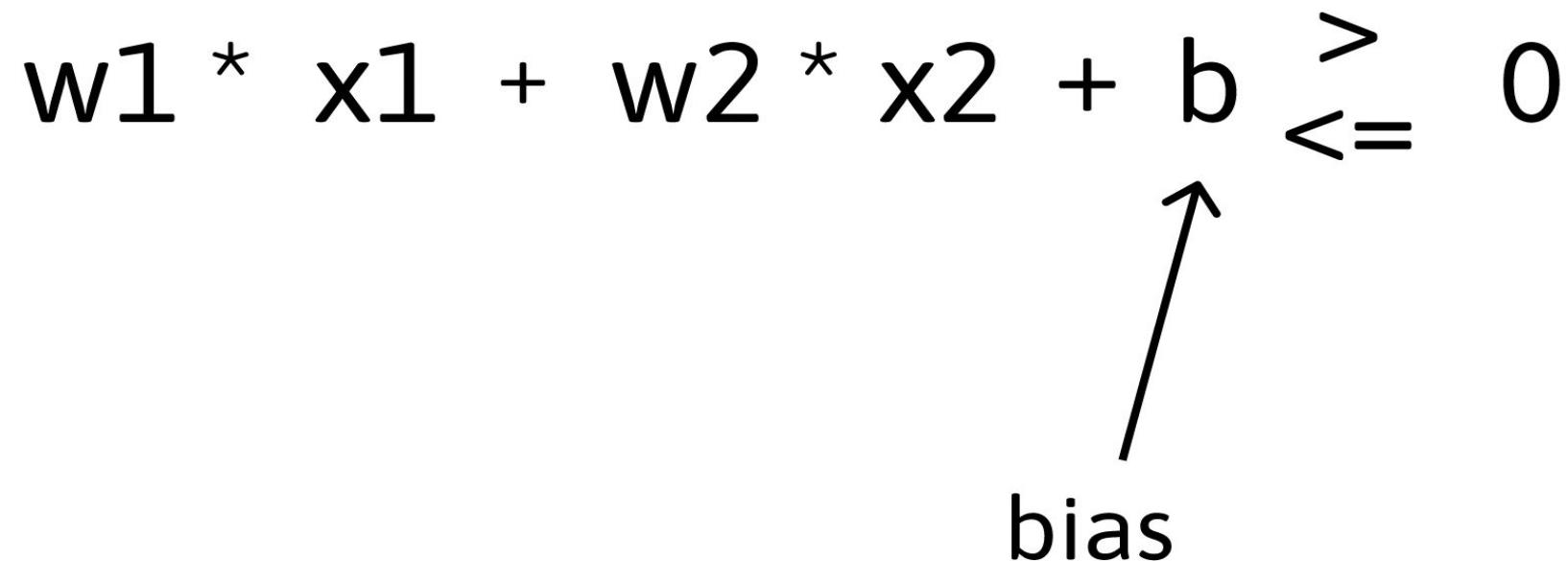


Price house prediction

Decide whether is an expensive house or not?

$$w_1 * x_1 + w_2 * x_2 + b \begin{matrix} > \\ \leq \end{matrix} 0$$

bias



Price house prediction

Solving the problem

Training set

(20, 10) 1

(10, 10) 1

(5, 5) 0

Test set

(15, 7)

Price house prediction

Solving the problem

Training set

(20, 10) 1
(10, 10) 1
(5, 5) 0



Determine w_1 , w_2 and b such that

$$w_1 * 20 + w_2 * 10 + b > 0 \text{ and}$$
$$w_1 * 10 + w_2 * 10 + b > 0 \text{ and}$$
$$w_1 * 5 + w_2 * 5 + b \leq 0$$

at the same time

Test set

(15, 7)

Price house prediction

Solving the problem

Training set

(20, 10) 1
(10, 10) 1
(5, 5) 0



Determine w_1 , w_2 and b such that

$$w_1 * 20 + w_2 * 10 + b > 0 \text{ and}$$
$$w_1 * 10 + w_2 * 10 + b > 0 \text{ and}$$
$$w_1 * 5 + w_2 * 5 + b \leq 0$$

at the same time

Test set

(15, 7)



See whether

$$w_1 * 15 + w_2 * 7 + b \stackrel{>}{\stackrel{<=}{}} 0$$

Price house prediction

Solving the problem

Training set

(20, 10) 1
(10, 10) 1
(5, 5) 0



Determine w_1 , w_2 and b such that
 $w_1 * 20 + w_2 * 10 + b > 0$ and
 $w_1 * 10 + w_2 * 10 + b > 0$ and
 $w_1 * 5 + w_2 * 5 + b \leq 0$
at the same time

Test set

(15, 7)



See whether
 $w_1 * 15 + w_2 * 7 + b \begin{matrix} > \\ \leq \end{matrix} 0$

Price house prediction

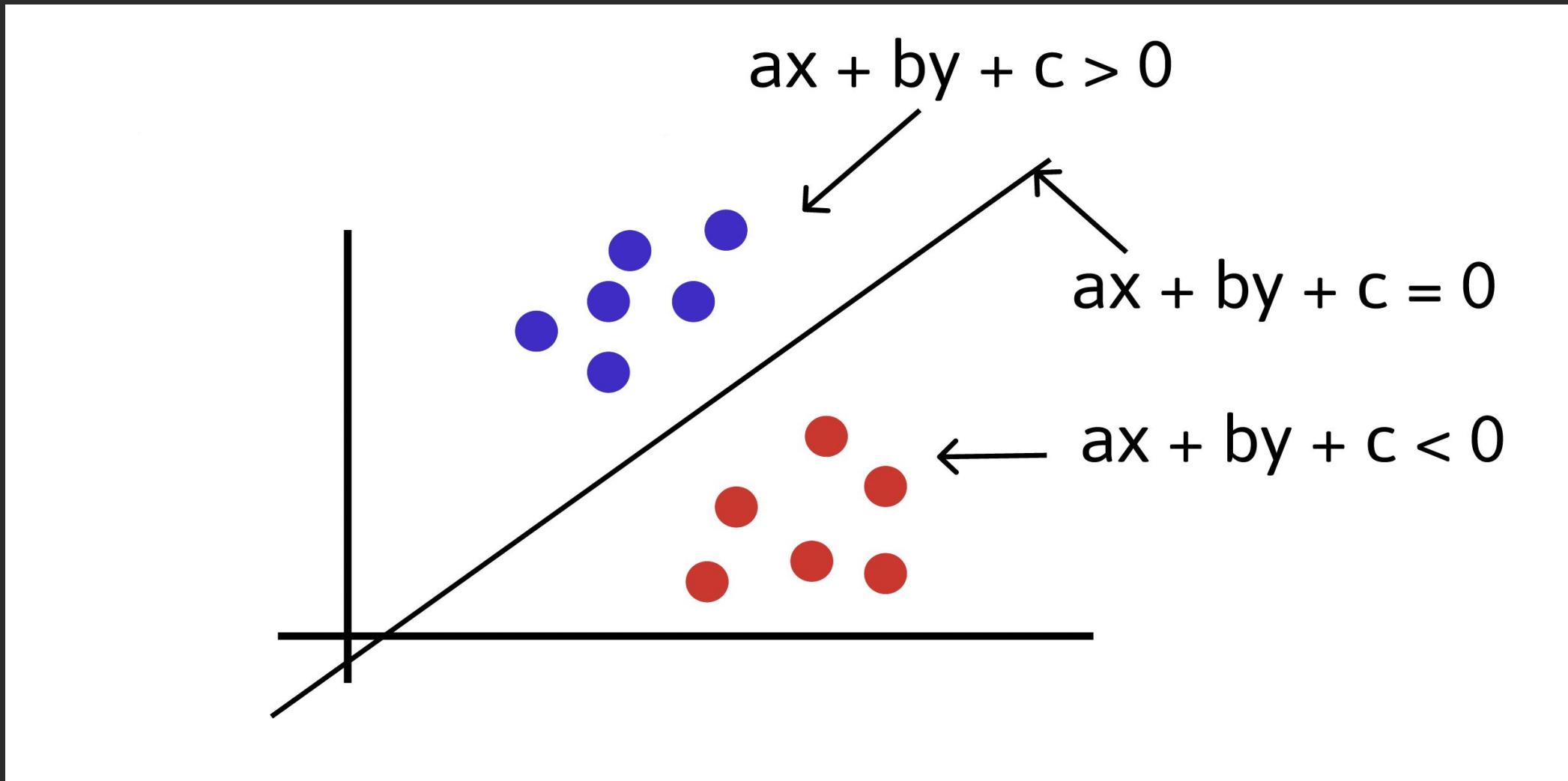
Solving the problem

$$w_1 * x_1 + w_2 * x_2 + b \geq 0$$
$$w_1 * x_1 + w_2 * x_2 + b \leq 0$$
$$ax + by + c = 0$$

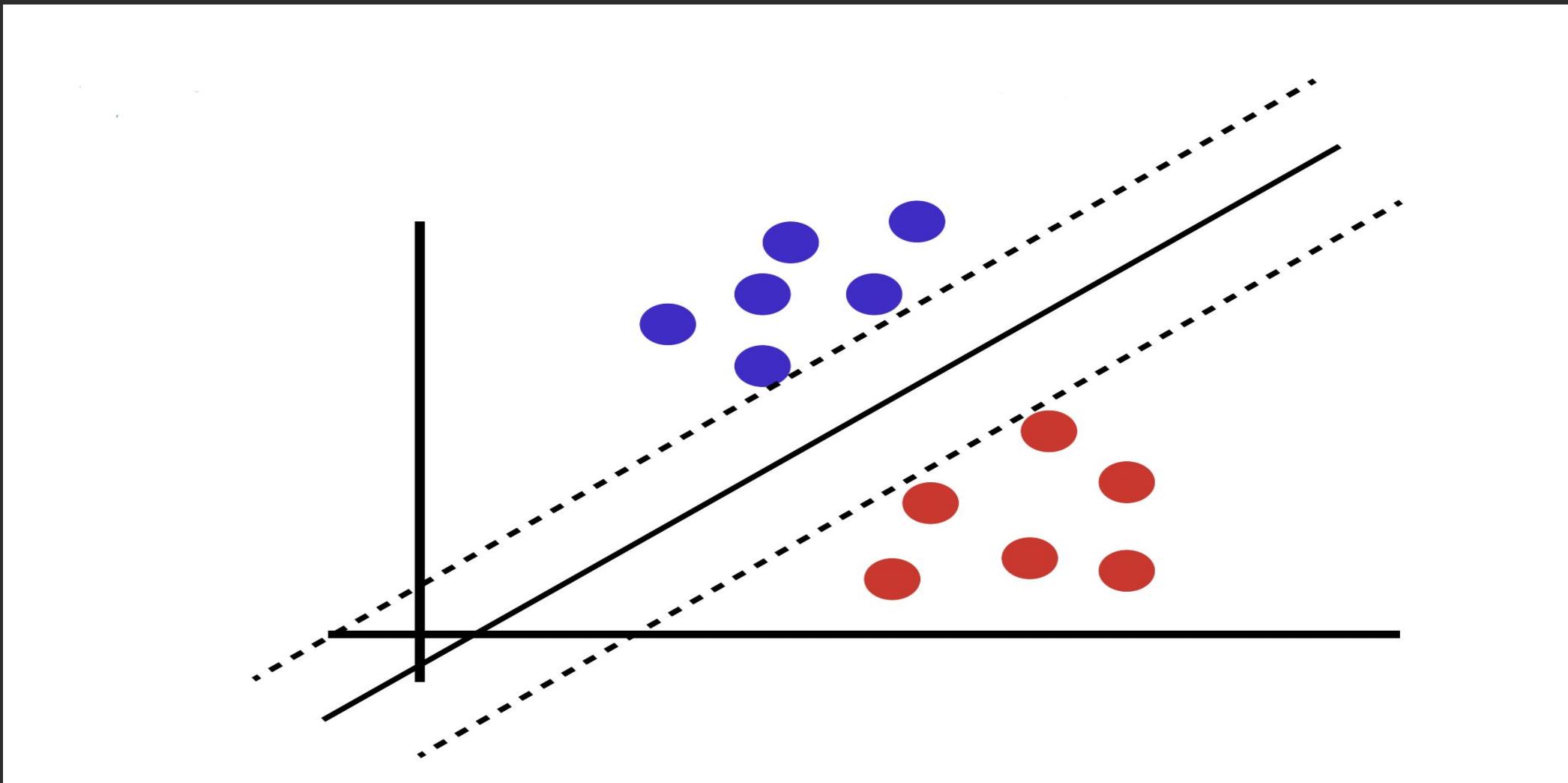
The diagram illustrates the relationship between a linear inequality and a linear equation. It shows two equations above a third equation. Arrows point from the terms $w_1 * x_1$, $w_2 * x_2$, and b in the first two equations down to the corresponding terms ax , by , and c in the third equation, indicating that these terms are equal in both the inequality and the equation.

Price house prediction

Solving the problem



Support Vector Machines



Support Vector Machines

Training set

(20, 10) 1

(10, 10) 1

(5, 5) 0



Determine w_1, w_2 and b such that

$$w_1 * 20 + w_2 * 10 + b > 0 \text{ and}$$

$$w_1 * 10 + w_2 * 10 + b > 0 \text{ and}$$

$$w_1 * 5 + w_2 * 5 + b \leq 0$$

at the same time

While minimizing $\|w\|$, where $w = (w_1, w_2, b)$

Support Vector Machines

Training set

(20, 10) 1
(10, 10) 1
(5, 5) 0



Determine w_1 , w_2 and b such that

$$w_1 * 20 + w_2 * 10 + b \geq 1 \text{ and}$$

$$w_1 * 10 + w_2 * 10 + b \geq 1 \text{ and}$$

$$w_1 * 5 + w_2 * 5 + b \leq -1$$

at the same time

While minimizing $\|w\|$, where $w = (w_1, w_2, b)$

How do we know if our model is good? (Validation)

- *Split the available annotated data, e.g. 90% for training, 10% for validation*
- *Train the model on the training data and test the current configuration on the validation data*
- *Choose the model with the best performance on the validation data*

Tips for competition: You can also train the best performing model on all the data (train + val) before submitting

Q&A

Thank you!