# Data and Features, Our First ML Algo: LwP

CS771: Introduction to Machine Learning

# Announcements

- Please join on Piazza
  - Joining link already shared in lecture 1 slides
  - If you don't join by Aug 10, there will be some penalty
  - If you have some issue joining Piazza, let us know and we can add you

- Some bonus marks will be reserved for (constructive) Piazza participation
  - Insightful questions and helpful answers/discussions from students
  - Amount of bonus at instructor's discretion (e.g., if you are falling below the threshold of a grade, you may get a "bump up")

- Project groups should be declared by Aug 11
  - Will share a Google form to enter details

# Plan today

- Feature extraction from raw data
  - How to convert raw inputs into "features" that our ML algos can understand/use?

    **Images as inputs**

    **Text docs as inputs**

  - How to transform some given features to make them more useful?

- Our first machine learning algorithm
  - Learning with Prototypes (LwP)

# Supervised Learning: Regression and Classification

- Regression is when the output is real-valued, e.g.,

Other types of supervised learning problems also exist, such as ranking, which we will study later

Input =     Output = Weight of the dog

- Classification is when the output is discrete, e.g.,

Input =     Output = Is it a dog (1) or cat (0) ?

Binary classification

Input =     Output = Breed (one of $K > 2$ possibilities)

Label is a "one-hot" vector of length $K$

Multi-class classification

Label is a binary vector of length $L$

$L > 1$ labels, each binary

Input =     Output = Dog (1/0)? Fluffy (1/0)? Sitting (1/0)? Wearing collar (1/0)?

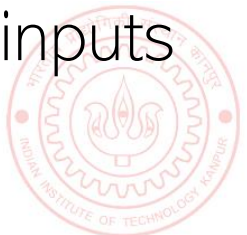Also called "tagging"    Multi-label classification

$L = 4$ in this case

# Data and Features

Features represent semantics of the inputs. Being able to extract good features is key to the success of ML algos
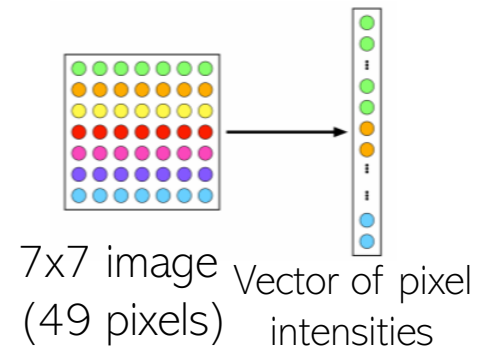
- ML algos require a numeric feature representation of the inputs

- Features can be obtained using one of the two approaches
    - Approach 1: Extracting/constructing features <u>manually</u> from raw inputs
    - Approach 2: <u>Learning</u> the features from raw inputs

- Approach 1 is what we will assume/focus on primarily for now
    - For the first few lectures, we will assume that someone has already given us the features

- Approach 2 is what is followed in Deep Learning based models and other ML models that learn a "code" (some optimal feature representation) for the inputs

- Approach 1 is not as powerful as Approach 2 but still used widely
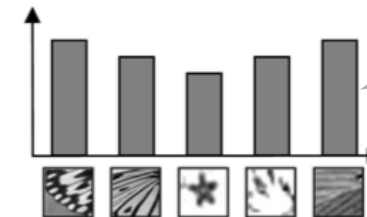
# Example: Feature Extraction for Image Data

- A very simple feature extraction approach for image data is <span style="color:red">flattening</span>

7x7 image (49 pixels)

Vector of pixel intensities

Flattening and histogram based methods destroy the spatial information in the image but often still work reasonably well

- <span style="color:red">Histogram</span> of visual patterns is another popular feature extr. method for images

Bar heights in the histogram denote how the <span style="color:red">frequency of occurrence</span> of each of the patterns in the given image (this vector of frequencies can be used as the extracted feature vector for this image)

These patterns can also be "discovered" by some feature learning algorithms (will see later)

Suppose these are typical patterns in the images in the dataset (some "domain expert" may have told us)
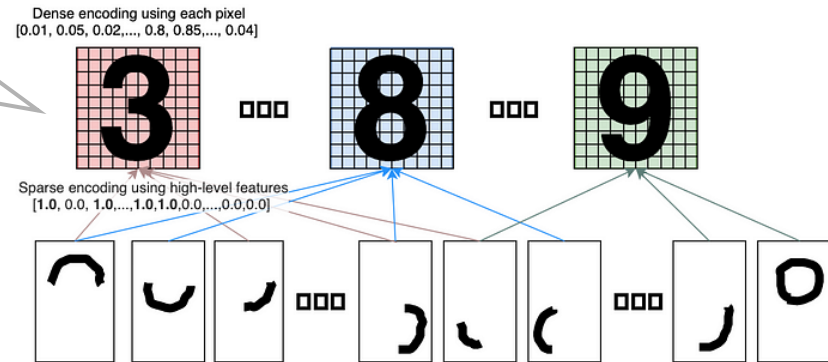
- Many other manual feature extraction techniques developed in computer vision and image processing communities (SIFT, HoG, and others)

# Example: Feature Extraction for Image Data

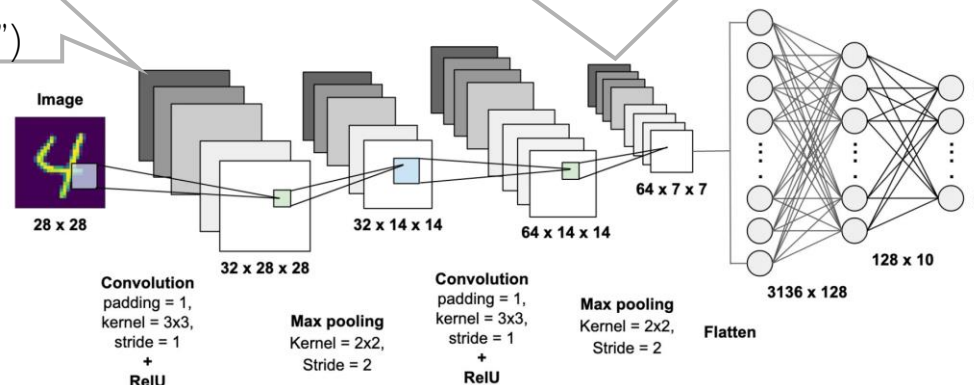- (Sparse) coding methods can also be used to <u>learn</u> good features

Each input represented using a binary feature vector denoting presence/absence of a large number of patterns (handwriting strokes)

Dense encoding using each pixel
[0.01, 0.05, 0.02,..., 0.8, 0.85,..., 0.04]

Sparse encoding using high-level features
[1.0, 0.0, 1.0,...,1.0,1.0,0.0,...,0.0,0.0]

- Deep Learning methods offer a very powerful way to <u>learn</u> good features

Last layer of features (the final supervised learning model uses these features)

Several layers of feature representations learned for each input (that's why it's called "deep")

Image
28 x 28

32 x 28 x 28

Convolution
padding = 1,
kernel = 3x3,
stride = 1
+
ReLU

Max pooling
Kernel = 2x2,
Stride = 2

32 x 14 x 14

Convolution
padding = 1,
kernel = 3x3,
stride = 1
+
ReLU

64 x 14 x 14

Max pooling
Kernel = 2x2,
Stride = 2

64 x 7 x 7

Flatten

3136 x 128

128 x 10

0
1
:
9

CS771: Intro to ML

# Example: Feature Extraction for Text Data

- Consider some text data consisting of the following sentences:
  - John likes to watch movies
  - Mary likes movies too
  - John also likes football

> BoW is just one of the many ways of doing feature extraction for text data. Not the most optimal one, and has various flaws (can you think of some?), but often works reasonably well

- Want to construct a feature representation for these sentences

- Here is a "bag-of-words" (BoW) feature representation of these sentences

Our "vocabulary" of 9 words

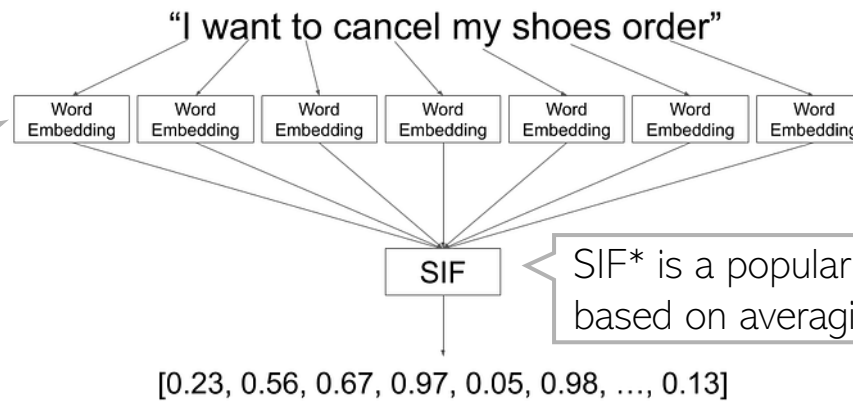| | John | likes | to | watch | movies | Mary | too | also | football |
|---|---|---|---|---|---|---|---|---|---|
| Sentence 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| Sentence 2 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| Sentence 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

- Each sentence is now represented as a binary vector (each feature is a binary value, denoting presence or absence of a word). BoW is also called "unigram" rep.
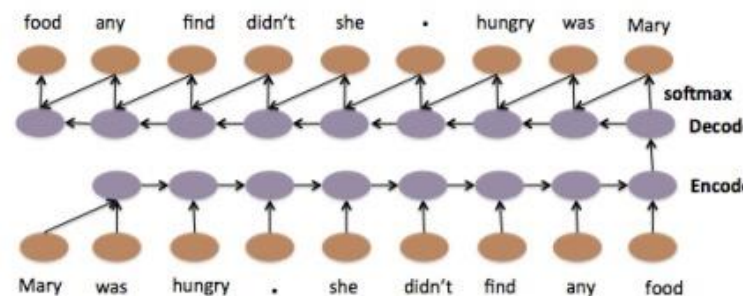
# Example: Feature Extraction for Text Data

- For text data (like sentences or document), we can do feature extraction using other simple methods that also use explicit semantics of words

For each word, there exists an "embedding" (a vector) such that semantically similar words have similar embeddings

"I want to cancel my shoes order"

| Word Embedding | Word Embedding | Word Embedding | Word Embedding | Word Embedding | Word Embedding | Word Embedding |

SIF

SIF* is a popular feature extraction method based on averaging of word embeddings

[0.23, 0.56, 0.67, 0.97, 0.05, 0.98, …, 0.13]

- Deep learning methods for sequence data (e.g., recurrent neural networks and transformers) are more powerful for feature extraction for text data (more on this later)

food any find didn't she . hungry was Mary

softmax

Decode

Encode

Mary was hungry . she didn't find any food

The state representation at the end of this sequence can be used as a feature for this sentence

*A Simple but Tough-to-Beat Baseline for Sentence Embeddings (Arora et al, 2017)

CS771: Intro to ML

# Don't forget to apply common sense with features!

- Consider a problem of predicting some air-pollutant's (e.g., PM2.5) concentration at various locations and at various times

- Here, in addition to other features, we also have features such as
  - Lat-long of each location
  - Time-stamp (e.g., hour of the day)

- Using raw values of these features may not be helpful (could even hurt) since
  - Earth is round ☺
  - Hours of the day have cyclical relationships (hour 23 is closer to hour 1 than hour 18)

- Thus we may need to transform such features using suitable techniques
  - Cyclical embeddings and other specific techniques for feature transformation can be used

- Caution: Don't leave deep learning to take care of all the features. Do apply some common sense as well

# Feature Learning = Distance Function Learning?

- It seems there are two ways to learn effectively from data
  - Extract (or rather "learn") features from the data + use standard distance function

Assume $f$ represents is our feature extraction function

$$d(\boldsymbol{a}, \boldsymbol{b}) = \sqrt{(f(\boldsymbol{a}) - f(\boldsymbol{b}))^\top (f(\boldsymbol{a}) - f(\boldsymbol{b}))}$$

  - Use a good distance/similarity function of the "basic" features, e.g.,

$$d_w(\boldsymbol{a}, \boldsymbol{b}) = \sqrt{(\boldsymbol{a} - \boldsymbol{b})^\top \mathbf{W} (\boldsymbol{a} - \boldsymbol{b})}$$

$$k(\boldsymbol{a}, \boldsymbol{b}) = \exp(-\gamma \|\boldsymbol{a} - \boldsymbol{b}\|^2)$$

- Both approaches, at some level, can be seen doing the same thing, just in different ways
- Any feature learning method corresponds to learning some distance/similarity function
  - .. and vice-versa

# Feature Selection

- Not all the extracted features may be relevant for learning the model (some may even confuse the learner)

- Feature selection (a step after feature extraction) can be used to identify the features that matter, and discard the others, for more effective learning

~~Age~~
~~Gender~~
Height
Weight
~~Eye color~~

⟹ Body-mass index (BMI)

Calculating BMI from this data doesn't require ML but this simple example is just to illustrate the idea of feature selection ☺

- Many techniques exist – some based on intuition, some based on algorithmic principles (will visit feature selection later)

- More common in supervised learning but can also be done for unsup. Learning
- Many ML algorithms can automatically take care of feature extraction/selection
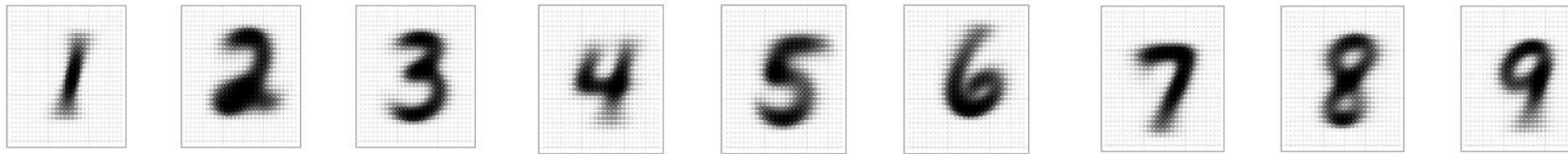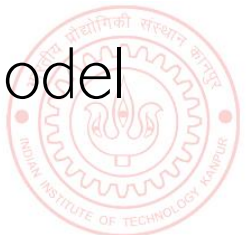  - More on this later

# Our First ML algorithm

# Classification via Learning with Prototypes (LwP)

- Basic idea: Represent each class by a "prototype" vector

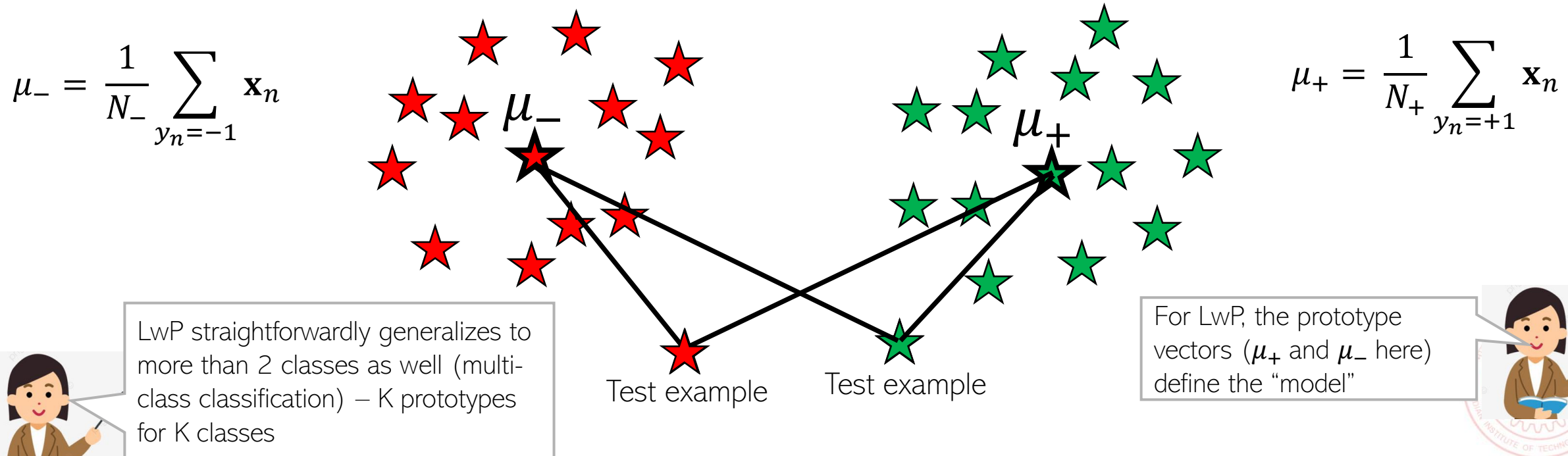- Class Prototype: The "mean" or "average" of inputs from that class



Averages (prototypes) of each of the handwritten digits 1-9

- Predict label of each test input based on its distances from the class prototypes
    - Predicted label will be the class that is the closest to the test input

- How we compute distances can have an effect on the accuracy of this model (may need to try Euclidean, weight Euclidean, or something else)

# Learning with Prototypes (LwP): An Illustration

- Suppose the task is binary classification (two classes assumed pos and neg)

- Training data: $N$ labelled examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \{-1, +1\}$
  - Assume $N_+$ example from positive class, $N_-$ examples from negative class
  - Assume green is positive and red is negative

$$\mu_- = \frac{1}{N_-} \sum_{y_n=-1} \mathbf{x}_n \qquad \mu_+ = \frac{1}{N_+} \sum_{y_n=+1} \mathbf{x}_n$$



$\mu_-$

$\mu_+$

Test example    Test example

LwP straightforwardly generalizes to more than 2 classes as well (multi-class classification) – K prototypes for K classes

For LwP, the prototype vectors ($\boldsymbol{\mu}_+$ and $\boldsymbol{\mu}_-$ here) define the "model"
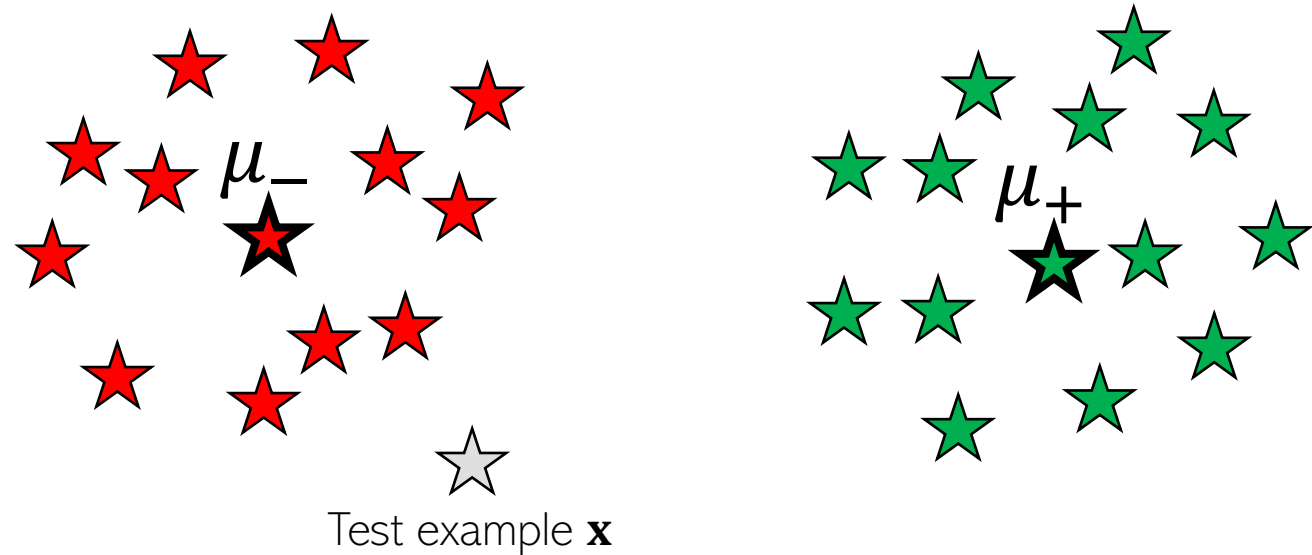
# LwP: The Prediction Rule, Mathematically

- What does the prediction rule for LwP look like mathematically?

- Assume we are using Euclidean distances here

$$||\boldsymbol{\mu}_- - \mathbf{x}||^2 = ||\boldsymbol{\mu}_-||^2 + ||\mathbf{x}||^2 - 2\langle \boldsymbol{\mu}_-, \mathbf{x}\rangle$$

$$||\boldsymbol{\mu}_+ - \mathbf{x}||^2 = ||\boldsymbol{\mu}_+||^2 + ||\mathbf{x}||^2 - 2\langle \boldsymbol{\mu}_+, \mathbf{x}\rangle$$

$\mu_-$

$\mu_+$

Test example $\mathbf{x}$

**Prediction Rule:** Predict label as +1 if $f(\mathbf{x}) = ||\boldsymbol{\mu}_- - \mathbf{x}||^2 - ||\boldsymbol{\mu}_+ - \mathbf{x}||^2 > 0$ otherwise -1

# LwP: The Prediction Rule, Mathematically

- Let's expand the prediction rule expression a bit more

$$\begin{aligned}
f(\mathbf{x}) &= ||\boldsymbol{\mu}_- - \mathbf{x}||^2 - ||\boldsymbol{\mu}_+ - \mathbf{x}||^2 \\
&= ||\boldsymbol{\mu}_-||^2 + ||\mathbf{x}||^2 - 2\langle\boldsymbol{\mu}_-, \mathbf{x}\rangle - ||\boldsymbol{\mu}_+||^2 - ||\mathbf{x}||^2 + 2\langle\boldsymbol{\mu}_+, \mathbf{x}\rangle \\
&= 2\langle\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-, \mathbf{x}\rangle + ||\boldsymbol{\mu}_-||^2 - ||\boldsymbol{\mu}_+||^2 \\
&= \langle\mathbf{w}, \mathbf{x}\rangle + b
\end{aligned}$$

- Thus LwP with Euclidean distance is equivalent to a <span style="color:red">linear model</span> with
  - Weight vector $\mathbf{w} = 2(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-)$
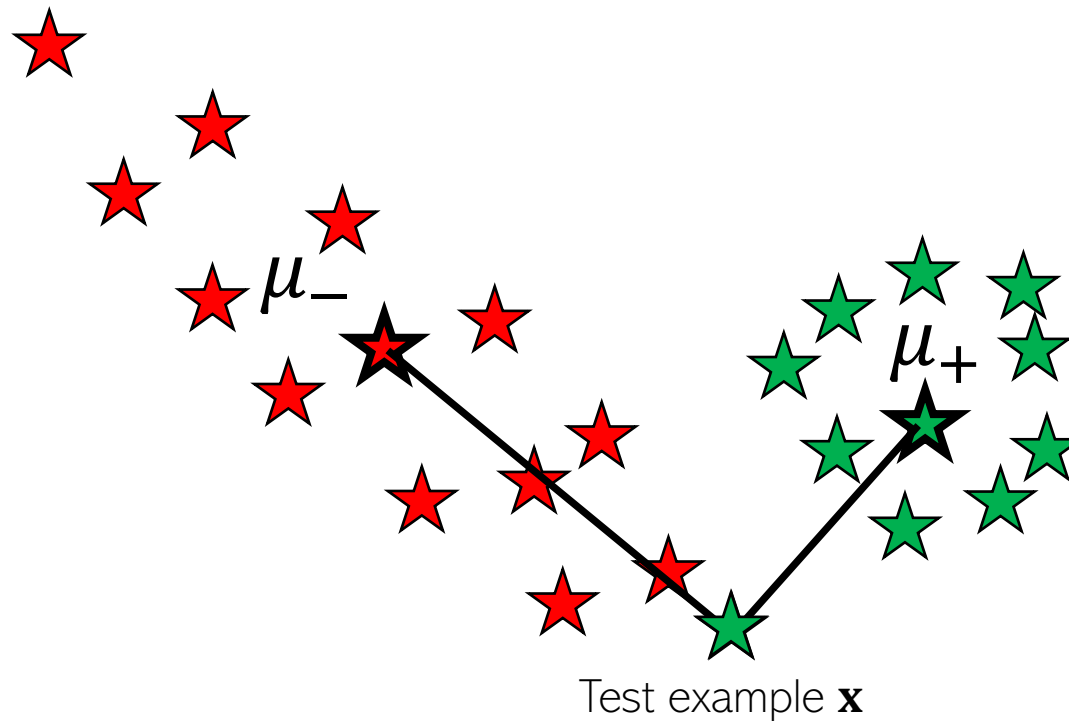  - Bias term $b = ||\boldsymbol{\mu}_-||^2 - ||\boldsymbol{\mu}_+||^2$

Will look at linear models more formally and in more detail later

- Prediction rule therefore is: Predict +1 if $\langle\mathbf{w}, \mathbf{x}\rangle + b > 0$, else predict -1

# LwP: Some Failure Cases

- Here is a case where LwP with Euclidean distance may not work well



Can use feature scaling or use Mahalanobis distance to handle such cases (will discuss this in the next lecture)

$\mu_-$

$\mu_+$

Test example **x**

- In general, if classes are not equisized and spherical, LwP with Euclidean distance will usually not work well. Can you think of how to fix this issue?

# LwP: Some Key Aspects

- Very simple, interpretable, and lightweight model
  - Just requires computing and storing the class prototype vectors

- Works with any number of classes (thus for multi-class classification as well)

- Can be generalized in various ways to improve it further, e.g.,
  - Modeling each class by a probability distribution rather than just a prototype vector
  - Using distances other than the standard Euclidean distance (e.g., Mahalanobis)

- With a learned distance function, can work very well even with very few examples from each class (used in some "few-shot learning" models nowadays – if interested, please refer to "Prototypical Networks for Few-shot Learning")