

Application: Few-shot classification

Classification with dynamically changing class labels. Do not need to decide on the set of classes during training time.

Example:

Test - 1

few = 2

Classes: {dog, panther, cat}

x^1 , dog

x^2 , dog

x^3 , panther

x^4 , panther

x^5 , cat

x^6 , cat

x^*

Test - 2

classes: {upma, pongal}

x^1 upma

x^2 upma

x^3 pongal

x^4 pongal

x^*

?

upma

pongal?

Learning distance functions

Many methods have been proposed.

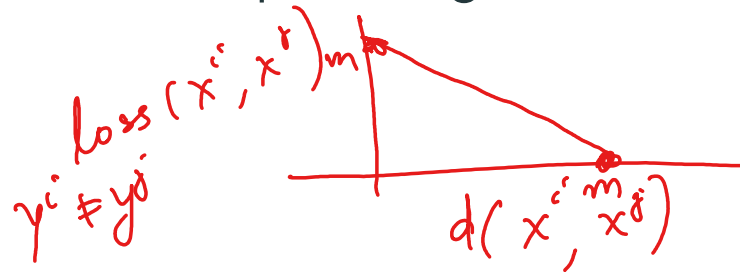
Assume we have supervision on pairs of objects that are similar and dissimilar.

Learn a distance function that brings together similar objects and keeps apart dissimilar objects.

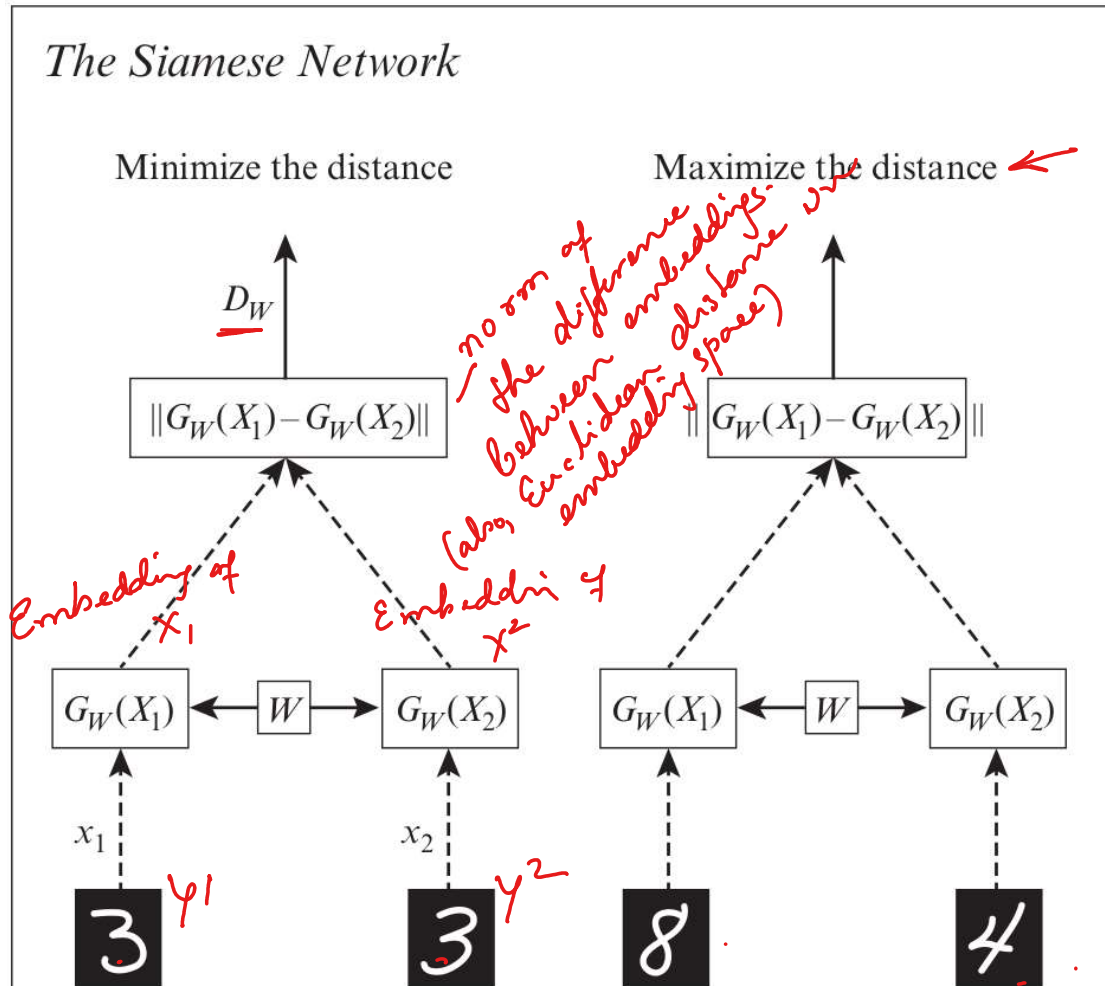
Default: $\min [(y_i = y_j) d(\mathbf{x}^i, \mathbf{x}^j) - \delta(y_i \neq y_j) d(\mathbf{x}^i, \mathbf{x}^j)]$ margin = hyper-parameter

$$\min L(\mathbf{x}^i, \mathbf{x}^j) = \delta(y_i = y_j) d(\mathbf{x}^i, \mathbf{x}^j) + \delta(y_i \neq y_j) \max(0, m - d(\mathbf{x}^i, \mathbf{x}^j))$$

Many different models proposed in traditional and deep learning methods for representing distances.



Siamese networks for learning distance functions



$$\max(0, m - D(G_W(x^1) - G_W(x^2)))$$

if $y^1 \neq y^2$

During training:

Use labeled data to learn W

During test-time:

Can support few shot classification with $\|G_W(x^i) - G_W(x^j)\|$ as distance function $d(x^i, x^j)$

Kernel Regression

Predict real-values based on similarity weighted values of training examples.

$$f(\mathbf{x}|D) = \sum_{i=1}^N \frac{K(\mathbf{x}, \mathbf{x}^i)}{\sum_{j=1}^N K(\mathbf{x}, \mathbf{x}_j)} y_i$$

Diagram illustrating the kernel regression formula:

A test point x^* is shown, with lines connecting it to training points x^1, x^2, \dots, x^N . The weights are proportional to the kernel values $k(x^*, x^i)$.

The training data is represented as a table:

x^1	y^1
x^2	y^2
\vdots	\vdots
x^N	y^N

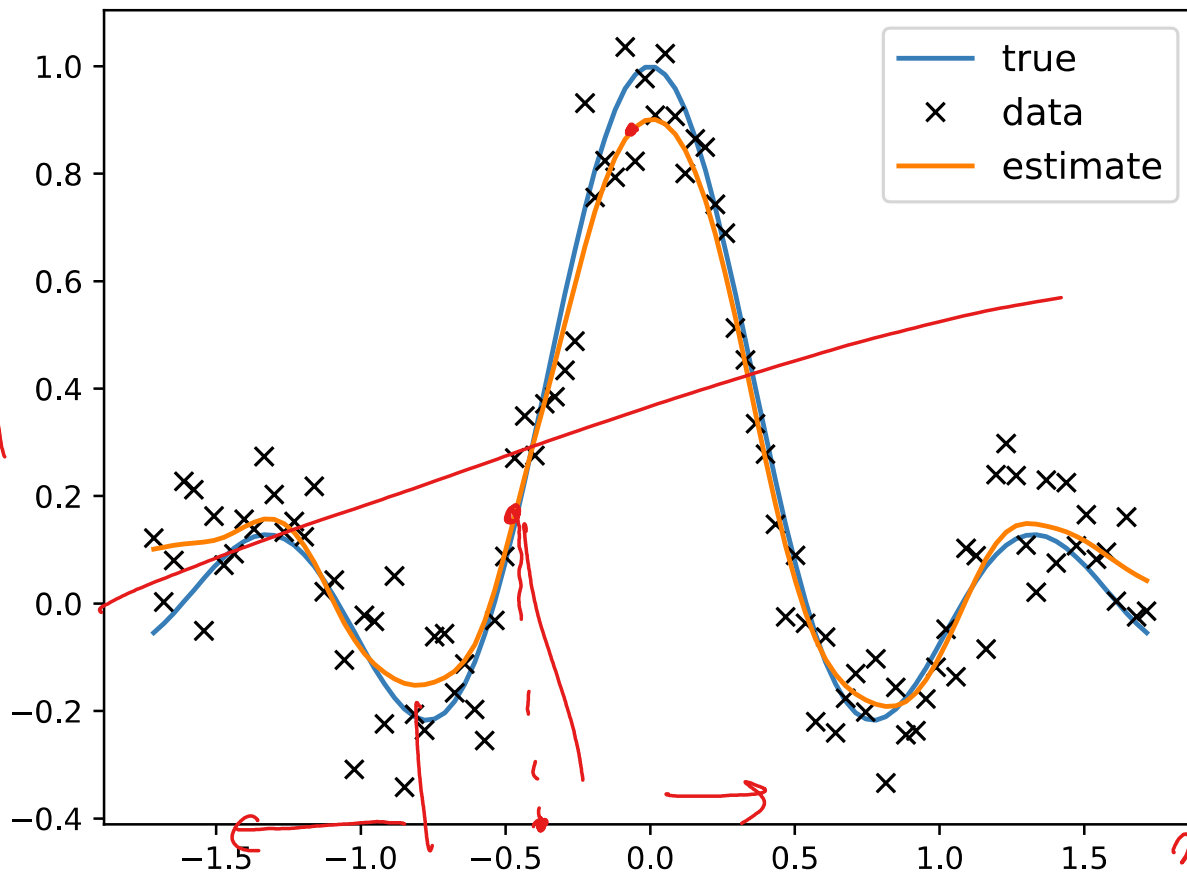
The final prediction is the weighted sum of the training values:

$$\hat{y} = \sum_{i=1}^N w_i y_i$$

The weights are calculated as:

$$w_i = \frac{k(x^*, x^i)}{\sum_{j=1}^N k(x^*, x^j)}$$

Kernel regression example



$d=1$
 $k(x^i, x^j) = e^{-\frac{r \|x^i - x^j\|^2}{2}}$
 $w_i = \frac{e^{-r \|x^i - x^*\|^2}}{\sum_k e^{-r \|x^k - x^*\|^2}}$
 if $r=0$, $\hat{y} = \frac{\sum_{i=1}^N y_i}{N}$
 if $r \rightarrow \infty$, $\hat{y} = y_{i^*}$
 $i^* = \text{nearest example to } x^*$

Mercer's Kernel function

The (Mercer's) kernel function $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is a symmetric function such that for any set of N points and any choice of numbers $c_i \in \mathbb{R}$ $k(x^i, x^j) = k(x^j, x^i)$

$$\sum_{i=1}^N \sum_{j=1}^N \underline{K(x^i, x^j)} \underline{c_i c_j} \geq 0$$

Equivalent to matrix of kernel values (Gram) matrix being positive semi-definite.

$$N \begin{bmatrix} k(x^i, x^j) \\ N \end{bmatrix}$$

Commonly Used Kernels

- Linear Kernel: $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$
- Polynomial Kernel: $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^d$ *↖ degree of polynomial.*
- Radial Basis Function (RBF) or Gaussian Kernel:
 $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$
- Sigmoid Kernel: $K(\mathbf{x}, \mathbf{y}) = \tanh(\alpha \mathbf{x}^T \mathbf{y} + c)$

Kernel embedding

Mercer's

For every valid Kernel, there exists an embedding function $\phi(\mathbf{x})$ such that

$$\underline{K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}') \text{ (Mercer's theorem)}}$$

(1cm) Example: Quadratic Kernel: $\underline{K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}')^2}$.

If original data is 2-d, then $\underline{K(\mathbf{x}, \mathbf{x}') = (x_1 x'_1 + x_2 x'_2)^2 = (x_1 x'_1)^2 + (x_2 x'_2)^2 + 2 x_1 x'_1 x_2 x'_2}$

$$\underline{\phi(\mathbf{x}) = [x_1^2, \sqrt{2} x_1 x_2, x_2^2]} \in \mathbb{R}^3$$

$$\phi(\mathbf{x}') = [x_1'^2, \sqrt{2} x_1' x_2', x_2'^2] \in \mathbb{R}^3$$

$$\underline{\phi(\mathbf{x}) \cdot \phi(\mathbf{x}') = x_1^2 \cdot x_1'^2 + (\sqrt{2} x_1 x_2)(\sqrt{2} x_1' x_2') + x_2^2 \cdot (x_2')^2}$$

$\stackrel{1}{=}$

$$= \underbrace{x_1^2 \cdot x_1'^2 + x_2^2 \cdot x_2'^2}_{\text{from above}} + (\sqrt{2} x_1 x_2) \cdot (\sqrt{2} x_1' x_2')$$

Applications in Machine Learning

- Support Vector Machines (SVM)
- Kernel Principal Component Analysis (PCA)
- Gaussian Processes
- Radial Basis Function Networks (RBFN)

Support Vector Machines

Non-probabilistic models for classification of the form:

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}, \mathbf{x}^i) + w_0$$

Learning fixes values of α_i . These are chosen in such a way that only a few i s have non-zero α_i . These are called support vectors.

Assume $y_i \in \{-1, +1\}$, binary classification model.

Predicted class label = $\text{sign}(f(\mathbf{x}))$