163

Lecture 13 - Reproducing Kernel Hilbert Spaces

OBJECTIVE: In this section, we derive one of the most fundamental theorems in machine learning. The theorem establishes a formal connection between kernels and dot products of feature vectors in a (possibly infinite) feature space.

Symmetric Positive Definite Matrices

If $\mathbf{A} = \mathbf{A}^{\mathbf{T}}$ (so \mathbf{A} has to be square!) then \mathbf{A} is said to be *symmetric*. If \mathbf{A} is symmetric, its eigenvalues are real.

A matrix \mathbf{A} is symmetric positive (semi)definite (SPD) if it is symmetric and

$$\mathbf{x}^{\mathrm{T}} \mathbf{A} \mathbf{x} \ge 0, \quad \forall \mathbf{x} \neq \mathbf{0}.$$

If **A** is SPD, its eigenvalues are positive.

CPSC-540: Machine Learning

Kernels as Measures of Similarity

Given the data $\{\mathbf{x}_{1:n}, y_{1:n}\}$, where $\mathbf{x}_i \in \mathcal{X}$ is an object and y its label, we measure the similarity between objects with the following function (called kernel):

$$k: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$$

The matrix of kernels **K** for a given dataset is called the **Gram matrix**. A kernel is SPD if the induced Gram matrix for any dataset is SPD.

Kernels allow us to map objects to function spaces:

$$\mathbf{x}_i \mapsto k(\cdot, \mathbf{x}_i)$$

*

165

Example: Kernel PCA

We can project items in high-dimensions to low dimensions, while preserving their topology. For example we can ensure that if two objects are far apart in high dimensions, then they will be far apart in low dimensions.

Initially, we have n data points in a d dimensional space in the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$. We map this data to the function (feature) space "matrix" $\mathbf{\Phi} \in \mathbb{R}^{n \times \infty}$ and consider its SVD: $\mathbf{\Phi} = \mathbf{U} \Sigma \mathbf{V}^T$. Of course, we can't compute the SVD of a a matrix with an infinite number of columns. However, if we use the kernel trick: $\mathbf{K} = \mathbf{\Phi} \mathbf{\Phi}^T$, then we have:

$$\mathbf{K} = \mathbf{\Phi} \mathbf{\Phi}^T = \mathbf{U} \Sigma^2 \mathbf{U}^T$$

So we can compute the eigenvalue decomposition of the nby n Gram matrix and then the k pricipal components are:

$$\mathbf{U}_k \Sigma_k = \mathbf{K} \mathbf{U}_k \Sigma_k^{-1}$$

Theorem 6 A function $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$, which is either continuous or has finite domain, can be decomposed as a dot product of features (in a Hilbert space):

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \boldsymbol{\phi}(\mathbf{x}_i), \boldsymbol{\phi}(\mathbf{x}_j) \rangle = \boldsymbol{\phi}(\mathbf{x}_i) \boldsymbol{\phi}(\mathbf{x}_j)^T$$

if and only if it is SPD.

Proof: First we prove the IF statement. That is, we begin by proving that if k admits the decomposition $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \boldsymbol{\phi}(\mathbf{x}_i), \boldsymbol{\phi}(\mathbf{x}_j) \rangle$, then k is SPD.

*

CPSC-540: Machine Learning	167
Next we prove the ONLY IF state. That is, if k is S	SPD, then
it can be expressed as the dot product of feature r	nappings.
The proof consits of the following three parts:	
1. Construct a vector (linear) space containing the	he images
of x under $\boldsymbol{\phi}$.	
 Define a dot product in this space and prove valid. 	that it is
3. Show that the dot product satisfies $k(\mathbf{x}_i, \mathbf{x}_j) =$	$\langle oldsymbol{\phi}(\mathbf{x}_i), oldsymbol{\phi}(\mathbf{x}_j)$
+ Step 1:	

 \star Step 2:

Step 1:

168

 \star Step 3:

169

CPSC-540: Machine Learning





The top-left plot shows the initial data \mathbf{x} , where only two points have been labelled. By running the active learning algorithm, the computer asks the user to enter the label for a point that could minimize the Bayes risk the most (the square in the top-right plot). The process is then repeated in the bottom-left plot. The final classification using only these four labels is perfect.

One can use Bayesian decision theory to extend the semisupervised learning approach to the active learning domain. In this approach, the class posterior probabilities $p(y_i^u | \mathbf{y}_l, \mathbf{x})$ are approximated with the estimates $0 \le y_i^u \le 1$, yielding the following expression for the posterior risk of the Bayes classifier:

CPSC-540: Machine Learning

$$R(\mathbf{y}_u) = \sum_{i=1}^{M} \min(y_i^u, 1 - y_i^u)$$

★ Proof:		

To decide which unlabelled point \mathbf{x}_j requires a label y_j , we first solve the linear system to obtain \mathbf{y}_u . Adding the new label y_j to the training set forces us to have to solve the linear system again. Fortunately, we can recompute the labels efficiently using the matrix inversion lemma.

Let the field obtained by adding the point (\mathbf{x}_j, y_j) be denoted by $\mathbf{y}_u^{+(\mathbf{x}_j, y_j)}$. Then, the posterior risk is:

$$R(\mathbf{y}_{u}^{+(\mathbf{x}_{j},y_{j})}) = \sum_{i=1}^{M-1} \min\left(y_{i}^{+(\mathbf{x}_{j},y_{j})}, 1 - y_{i}^{+(\mathbf{x}_{j},y_{j})}\right).$$

Of course, before querying the user, we do not know the label y_j , so we have to use the expected posterior risk:

$$\mathbb{E}[R(\mathbf{y}_{u}^{+(\mathbf{x}_{j},y_{j})})] = (1 - y_{i}^{u})R(\mathbf{y}_{u}^{+(\mathbf{x}_{j},0)}) + y_{i}^{u}R(\mathbf{y}_{u}^{+(\mathbf{x}_{j},1)})$$

After computing this expression, we pick the index j that minimizes it and ask the user to enter a label for \mathbf{x}_j . The pseudo-code follows.

173

Input: L, U and \mathbf{W}

WHILE more labels are required:

- Compute \mathbf{y}_u using (1).
- Find the best query j using the expected posterior risk.
- Query point \mathbf{x}_i and receive answer y_i .
- Add (\mathbf{x}_i, y_i) to L and remove \mathbf{x}_i from U.

CPSC-540: Machine Learning

The Learning Challenges

- Active Learning.
- Learning and Game Theory.
- Exciting new data.
- Function spaces and first order probabilistic inference.
- Solving complex probabilistic models in high dimensions.
- Combining all the bricks to build the cathedral.