

COMS 4771: Machine Learning

Homework 4, due April 19.

In this homework, we will study kernelized forms of ridge regression and PCA.

1 Kernel Ridge Regression

We derive kernel ridge regression in a couple steps.

Step 1. Consider solving the ridge regression problem with a training set

$$S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^p \times \mathbb{R}\}.$$

Let $\hat{\mathbf{w}}_\lambda$ be the ridge regression solution with regularization parameter λ , i.e.

$$\hat{\mathbf{w}}_\lambda := \arg \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2 + \lambda \|\mathbf{w}\|_2^2.$$

In class, we derived a closed form expression for $\hat{\mathbf{w}}_\lambda$ in terms of the design matrix \mathbf{X} and response vector \mathbf{y} , defined as

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_n^\top & - \end{bmatrix} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

This closed form expression turns out to be

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

We will derive an alternative form for $\hat{\mathbf{w}}_\lambda$. First, note that the above closed form expression comes from solving the optimality equation for $\hat{\mathbf{w}}_\lambda$, i.e.

$$(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_p) \hat{\mathbf{w}}_\lambda = \mathbf{X}^\top \mathbf{y}.$$

where \mathbf{I}_p is the identity matrix of order p . The above equation can be rewritten as

$$n\lambda \hat{\mathbf{w}}_\lambda = \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\mathbf{w}}_\lambda). \tag{1}$$

Now, define $\boldsymbol{\alpha} := \mathbf{y} - \mathbf{X} \hat{\mathbf{w}}_\lambda$.

Extra credit (weightage: 2%). Using equation (1) and the definition of α , show that α satisfies the following equation:

$$\frac{1}{n\lambda} \mathbf{X} \mathbf{X}^\top \alpha + \alpha = \mathbf{y}. \quad (2)$$

Extra credit (weightage: 1%). Solve equation (2) for α , and use your solution in equation (1) to get the following alternative expression for $\hat{\mathbf{w}}_\lambda$:

$$\hat{\mathbf{w}}_\lambda = \frac{1}{n\lambda} \mathbf{X}^\top \left(\frac{1}{n\lambda} \mathbf{X} \mathbf{X}^\top + \mathbf{I}_n \right)^{-1} \mathbf{y} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + n\lambda \mathbf{I}_n)^{-1} \mathbf{y}. \quad (3)$$

Step 2. Now suppose labeled points (\mathbf{x}, y) are drawn from $\mathcal{X} \times \mathbb{R}$, where \mathcal{X} is the feature space. Suppose you are also given a kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ which computes the inner product for a feature map $\phi : \mathcal{X} \rightarrow \mathbb{H}$ where \mathbb{H} is the reproducing kernel Hilbert space¹ for k . In other words,

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle.$$

Suppose you have collected a training set of n examples

$$S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathcal{X} \times \mathbb{R}\}.$$

Let $\hat{\mathbf{w}}_\lambda$ be the solution of the ridge regression problem with regularization parameter λ when data points are mapped to \mathbb{H} using the feature mapping ϕ . In other words, consider doing ridge regression using the transformed training set

$$S' = \{(\phi(\mathbf{x}_1), y_1), (\phi(\mathbf{x}_2), y_2), \dots, (\phi(\mathbf{x}_n), y_n) \in \mathbb{H} \times \mathbb{R}\}.$$

Mathematically, $\hat{\mathbf{w}}_\lambda$ can be defined as

$$\hat{\mathbf{w}}_\lambda := \arg \min_{\mathbf{w} \in \mathbb{H}} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle)^2 + \lambda \|\mathbf{w}\|_2^2.$$

Part 1 of assignment: (weightage: 4%) Given a test point $\mathbf{x} \in \mathcal{X}$, give the pseudocode for computing the prediction using the ridge regression solution, i.e. $\langle \hat{\mathbf{w}}_\lambda, \phi(\mathbf{x}) \rangle$. You may use equation (3) for this even if you haven't attempted to prove the equation in the extra credit question.

Hint: construct the design matrix \mathbf{X} for the transformed training set S' , and observe that the matrix $\mathbf{X} \mathbf{X}^\top$ can be computed using the kernel function k .

2 Kernel PCA

Next, we turn to kernel PCA. We use the connection between PCA and SVD to derive a kernelized form of PCA.

Let \mathbf{X} be the design matrix, and let $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$ be the SVD of \mathbf{X} . Recall from class that the following facts:

1. The columns of \mathbf{U} are the eigenvectors of $\mathbf{X} \mathbf{X}^\top$, with eigenvalues given by the squares of the diagonal entries in \mathbf{S} .

¹If the jargon “reproducing kernel Hilbert space” alarms you, just think of \mathbb{H} as \mathbb{R}^D for some $D \gg p$.

2. The rank ℓ PCA representation of the training set is given by the columns of

$$\mathbf{S}_\ell^{-1} \mathbf{U}_\ell^\top \mathbf{X} \mathbf{X}^\top.$$

Extra credit (weightage: 3%). Now let \mathbf{X} be the design matrix for the transformed training set S' . Explain how you can compute \mathbf{S}_ℓ and \mathbf{U}_ℓ using the kernel function k without having to explicitly compute the mapping ϕ for any training points.

Part 2 of assignment: (weightage: 2%) Suppose you have computed \mathbf{S}_ℓ and \mathbf{U}_ℓ . Specify the dimensions of these two matrices. Explain how you can compute the rank ℓ PCA of the transformed training set using these matrices and the kernel function k without having to explicitly compute the mapping ϕ for any training points.