

Recitation 6: Kernel PCA, Ridge Regression

1 PCA

Recall in principal components analysis, we are interested in the following maximization problem. Given data: $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ that has been centered (each index of the data vectors has mean 0), the first principal component is the unit vector \mathbf{w} whose projection onto the data is maximized:

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_i^m (\mathbf{x}_i^\top \mathbf{w})^2$$

As we've seen in lecture, this is an eigenvalue problem in disguise. Let $\mathbf{X} = \begin{bmatrix} | & | & \dots & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_m \\ | & | & \dots & | \end{bmatrix}$ and rewrite the maximization problem in the form of an inner product of this data matrix.

$$\begin{aligned} \mathbf{w}_1 &= \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \|\mathbf{X}^\top \mathbf{w}\|^2 \\ &= \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \langle \mathbf{X}^\top \mathbf{w}, \mathbf{X}^\top \mathbf{w} \rangle \\ &= \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} \end{aligned}$$

This is exactly the Rayleigh Quotient that we saw in homework 1. So the \mathbf{w} that maximizes $\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}$ is the eigenvector with largest eigenvalue of $\mathbf{X}^\top \mathbf{X}$.

1.1 Kernel PCA

Now instead suppose our data comes from some set \mathcal{X} . We have some positive semidefinite kernel on this set $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, with the corresponding induced feature mapping $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$. We're interested in the same sort of maximization above but in the feature space given by ϕ .

Let $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}$ and $\Phi = \begin{bmatrix} | & | & \dots & | \\ \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_m) \\ | & | & \dots & | \end{bmatrix}$.

$$\begin{aligned} \mathbf{v}_1 &= \arg \max_{\|\mathbf{v}\|=1} \sum_{i=1}^m \langle \phi(\mathbf{x}_i), \mathbf{v} \rangle^2 \\ &= \arg \max_{\|\mathbf{v}\|=1} \|\Phi^\top \mathbf{v}\|^2 \\ &= \arg \max_{\|\mathbf{v}\|=1} \mathbf{v}^\top \Phi \Phi^\top \mathbf{v} \end{aligned}$$

So \mathbf{v}_1 is the top eigenvector of $\Phi\Phi^\top$. We can say a bit more about the top principal component vector \mathbf{v}_1 if we consider the form it must take with respect to the feature mappings of the input data.

Lemma 1 *The \mathbf{v}_1 that maximizes this sum of square projections onto the $\phi(\mathbf{x}_i)$'s will be a linear combination of the $\phi(\mathbf{x}_i)$'s:*

$$\mathbf{v}_1 = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i)$$

for some set of $\alpha_i \in \mathbb{R}$. This can be equivalently written as: $\mathbf{v}_1 = \Phi\boldsymbol{\alpha}$, for $\boldsymbol{\alpha} \in \mathbb{R}^m$.

Proof: Suppose \mathbf{v} is an eigenvector of $\Phi\Phi^\top = \sum_{i=1}^m \phi(\mathbf{x}_i)\phi(\mathbf{x}_i)^\top$:

$$\begin{aligned} \Phi\Phi^\top \mathbf{v} &= \sum_{i=1}^n \phi(\mathbf{x}_i)\phi(\mathbf{x}_i)^\top \mathbf{v} \\ \lambda \mathbf{v} &= \sum_{i=1}^n \phi(\mathbf{x}_i) \underbrace{\phi(\mathbf{x}_i)^\top \mathbf{v}}_{\in \mathbb{R}} \\ \mathbf{v} &= \frac{1}{\lambda} \sum_{i=1}^n (\phi(\mathbf{x}_i)^\top \mathbf{v}) \phi(\mathbf{x}_i) \end{aligned}$$

Let $\alpha_i = \frac{\phi(\mathbf{x}_i)^\top \mathbf{v}}{\lambda}$ and we have our desired result. ■

Now let's use this explicit form $\mathbf{v}_1 = \Phi\boldsymbol{\alpha}$ in our computations and see what pops out:

$$\begin{aligned} \Phi\Phi^\top \mathbf{v}_1 &= \Phi\Phi^\top \Phi\boldsymbol{\alpha} \\ \lambda \mathbf{v}_1 &= \Phi \mathbf{K} \boldsymbol{\alpha} \end{aligned}$$

where $\mathbf{K} = \Phi^\top \Phi \in \mathbb{R}^{m \times m}$ is the Gram matrix of our input data: $\mathbf{K}_{i,j} = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$.

Hit both sides by Φ^\top :

$$\begin{aligned} \Phi^\top \lambda \mathbf{v}_1 &= \Phi^\top \Phi \mathbf{K} \boldsymbol{\alpha} \\ \lambda \Phi^\top \Phi \boldsymbol{\alpha} &= \Phi^\top \Phi \mathbf{K} \boldsymbol{\alpha} \\ \lambda \mathbf{K} \boldsymbol{\alpha} &= \mathbf{K} \mathbf{K} \boldsymbol{\alpha} \end{aligned}$$

Eliminating one \mathbf{K} term from both sides gives us: $\mathbf{K}\boldsymbol{\alpha} = \lambda\boldsymbol{\alpha}$, which tells us that the coefficient vector $\boldsymbol{\alpha}$ of \mathbf{v}_1 , is in fact the top eigenvector of \mathbf{K} . Note that \mathbf{K} might not be full rank, but this will only be an issue for the zero-eigenvalued eigenvectors, which will not be a top principal component in the first place.

An equivalent way of getting to this result is by directly plugging in $\mathbf{v}_1 = \Phi\boldsymbol{\alpha}$ into our maximization problem:

$$\begin{aligned} \max_{\|\mathbf{v}\|=1} \mathbf{v}_1^\top \Phi\Phi^\top \mathbf{v}_1 &= \max_{\|\mathbf{v}\|=1} \boldsymbol{\alpha}^\top \Phi^\top \Phi \Phi^\top \Phi \boldsymbol{\alpha} \\ &= \max_{\|\mathbf{v}\|=1} \boldsymbol{\alpha}^\top \mathbf{K} \mathbf{K} \boldsymbol{\alpha} \end{aligned}$$

Thus, $\boldsymbol{\alpha}$ must be the top eigenvector of \mathbf{K}^2 . The eigenvectors of \mathbf{K}^2 are the same as the eigenvectors for \mathbf{K} . So computing the principal component of our data just boils down to finding the top eigenvector of the Gram matrix \mathbf{K} .

2 Kernel Ridge Regression

Using the same notation as in the previous section, suppose we have some base set \mathcal{X} , a positive semidefinite kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, and the corresponding Reproducing Kernel Hilbert Space \mathcal{H}_k induced by by our kernel k . Given $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}$, with corresponding target values $y_1, \dots, y_m \in \mathbb{R}$, we have the Gram matrix \mathbf{K} , with entries: $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$. We'd like to learn some function $f \in \mathcal{H}_k$ that fits our data subject to an additional regularization term:

$$\hat{f} = \arg \min_{f \in \mathcal{H}_k} \left[\underbrace{\sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2}_{\mathcal{R}[f]} + \lambda \|f\|_{\mathcal{H}_k} \right]$$

The Representer theorem tells us that \hat{f} must be of the form:

$$\hat{f}(\cdot) = \sum_{j=1}^m \alpha_j k(\cdot, \mathbf{x}_j)$$

Plug this form of \hat{f} into the minimization expression $\mathcal{R}[f]$. First let's consider the datafitting term of $\mathcal{R}[f]$:

$$\sum_{i=1}^m (\hat{f}(\mathbf{x}_i) - y_i)^2 = \sum_{i=1}^m \left(\sum_{j=1}^m \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) - y_i \right)^2 \quad (1)$$

$$= \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|^2 \quad (2)$$

where $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$, $\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix}$. The inner sum: $\sum_{j=1}^m \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$ is the same as taking the dot product of the i th row vector of \mathbf{K} with $\boldsymbol{\alpha}$.

Now expanding the regularization term:

$$\lambda \|f\|_{\mathcal{H}_k} = \lambda \left\langle \sum_{i=1}^m \alpha_i k(\cdot, \mathbf{x}_i), \sum_{j=1}^m \alpha_j k(\cdot, \mathbf{x}_j) \right\rangle \quad (3)$$

$$= \lambda \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (4)$$

$$= \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \quad (5)$$

where we use the definition of the inner product of \mathcal{H}_k : $\langle k(\cdot, x), k(\cdot, y) \rangle = k(x, y)$ and the linearity of the inner product to go from eq(3) to eq(4). Plugging eq(2) and eq(5) back into $\mathcal{R}[f]$:

$$\begin{aligned} \mathcal{R}[f] &= \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \\ &= \boldsymbol{\alpha}^\top \mathbf{K}^\top \mathbf{K} \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^\top \mathbf{K}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \end{aligned}$$

Taking the gradient of $\mathcal{R}[f]$ with respect to α and setting it to $\mathbf{0}$:

$$\begin{aligned}
 \nabla_{\alpha} \mathcal{R}[f] &= 2\mathbf{K}^{\top} \mathbf{K} \alpha - 2\mathbf{K}^{\top} \mathbf{y} + 2\lambda \mathbf{K} \alpha \\
 \mathbf{0} &= 2\mathbf{K}^{\top} \mathbf{K} \alpha - 2\mathbf{K}^{\top} \mathbf{y} + 2\lambda \mathbf{K} \alpha \\
 \mathbf{K}^2 \alpha + \lambda \mathbf{K} \alpha &= \mathbf{K} \mathbf{y} \\
 (\mathbf{K}^2 + \lambda \mathbf{K}) \alpha &= \mathbf{K} \mathbf{y} \\
 \alpha &= (\mathbf{K}^2 + \lambda \mathbf{K})^{-1} \mathbf{K} \mathbf{y} \\
 &= (\mathbf{K}(\mathbf{K} + \lambda \mathbf{I}))^{-1} \mathbf{K} \mathbf{y} \\
 &= (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^{-1} \mathbf{K} \mathbf{y} \\
 &= (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}
 \end{aligned}$$

\mathbf{K} is symmetric so we can replace \mathbf{K}^{\top} with \mathbf{K} in the equations above. Evaluating \hat{f} on some new point $\mathbf{z} \in \mathcal{X}$ then boils down to computing the inner product between $\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$ with the vector of kernel

evaluations: $\mathbf{k}_{\mathbf{z}} = \begin{bmatrix} k(\mathbf{z}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{z}, \mathbf{x}_m) \end{bmatrix}$

$$\hat{f}(\mathbf{z}) = \sum_{i=1}^m \alpha_i k(\mathbf{z}, \mathbf{x}_i) = \alpha^{\top} \mathbf{k}_{\mathbf{z}}$$

Something to think about: how does the penalty term $\lambda \|f\|_{\mathcal{H}_k}$ affect the regularized risk minimization problem? How does our solution, \hat{f} , change as we increase or decrease the parameter λ ?