

Machine Learning 1 Homework 13 Theory Part

01.02.2016

Kernel Ridge Regression

a) The cost function of Ridge Regression is defined as:

$$\mathcal{E}_{RR}(\mathbf{w}) = (\mathbf{y} - \mathbf{w}^T \mathbf{X})^2 + \lambda \|\mathbf{w}\|^2$$

\mathbf{X} : Matrix of input data points from training dataset, each column of it is a data point $\mathbf{x}_i \in \mathbb{R}^d$

\mathbf{y} : Vector of labels from training dataset

Take the derivative w.r.t \mathbf{w} yields:

$$\frac{\partial \mathcal{E}_{RR}(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}\mathbf{y} + 2\mathbf{X}\mathbf{X}^T \mathbf{w} + 2\lambda \mathbf{w}$$

We want to optimize this cost function, so set the derivative to 0, we get:

$$\begin{aligned}\mathbf{X}\mathbf{X}^T \mathbf{w} + \lambda \mathbf{w} &= \mathbf{X}\mathbf{y} \\ (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})\mathbf{w} &= \mathbf{X}\mathbf{y} \\ \mathbf{w} &= (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Now we get the solution for Ridge regression.

b) To kernelize this, we introduce kernel function: $k(\mathbf{x}_i, \mathbf{x}_j)$ which maps the innerproduct of two datapoints to another feature space.

The mapping can be denoted as:

$$\begin{aligned}\mathbf{x} &\rightarrow \Phi(\mathbf{x}) \\ \langle \mathbf{x}_i, \mathbf{x}_j \rangle &\rightarrow k(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)\end{aligned}$$

Thus we get the kernelized Ridge regression cost function:

$$\mathcal{E}_{RR}(\mathbf{w}) = (\mathbf{y} - \mathbf{w}^T \Phi(\mathbf{X}))^2 + \lambda \|\mathbf{w}\|^2$$

Take the derivative w.r.t \mathbf{w} and set it to 0, we get the solution:

$$\mathbf{w} = (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \Phi(\mathbf{X})\mathbf{y}$$

Combined with the linear assumption, we can get the estimated value given \mathbf{x} as input:

$$\begin{aligned}\hat{y} &= \Phi(\mathbf{x})^T \mathbf{w} \\ &= \Phi(\mathbf{x})^T (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \Phi(\mathbf{X})\mathbf{y} \\ &= \Phi(\mathbf{x})^T (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \Phi(\mathbf{X}) (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I}) (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \Phi(\mathbf{x})^T (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} (\Phi(\mathbf{X})\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I}\Phi(\mathbf{X})) (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \Phi(\mathbf{x})^T (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I}) \Phi(\mathbf{X}) (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \Phi(\mathbf{x})^T \Phi(\mathbf{X}) (\Phi(\mathbf{X})\Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \mathbf{k}^* (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \\ \mathbf{k}^* &= \Phi(\mathbf{x})^T \Phi(\mathbf{X}) \\ \mathbf{K} &= \Phi(\mathbf{X})\Phi(\mathbf{X})^T\end{aligned}$$

Since $(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$ can be calculated using training dataset, this part can be seen as coefficients α . \mathbf{k}^* is a vector which contains the mapped innerproduct of input data point \mathbf{x} with all the training data points. So this result can be written as the inner product of this vector \mathbf{k}^* with coefficients α . Thus we get the formula:

$$\hat{y} = \sum_{i=2}^n k(\mathbf{x}, \mathbf{x}_i) \alpha_i$$

c) Consider the new prime problem:

$$\begin{aligned} \min_{\xi, \mathbf{w}} \quad & \sum_{i=1}^n \xi_i^2 \\ \text{subject to} \quad & \xi_i = \mathbf{w}^T \mathbf{x}_i - y_i \\ & \|\mathbf{w}\|^2 \leq C \end{aligned}$$

Apply Lagrangian multiplier method:

$$\mathcal{L} = \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \lambda_i (\mathbf{w}^T \mathbf{x}_i - y_i - \xi_i) - \alpha (\|\mathbf{w}\|^2 - C)$$

Calculate derivatives w.r.t \mathbf{w} and ξ , we get:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \xi_i} &= 2\xi_i - \lambda_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \sum_{i=1}^n \lambda_i \mathbf{x}_i - 2\alpha \mathbf{w} = 0 \\ \Rightarrow \quad \xi_i &= \frac{1}{2} \lambda_i \\ \mathbf{w} &= \frac{1}{2\alpha} \sum_{i=1}^n \lambda_i \mathbf{x}_i \end{aligned}$$

insert the result to the primal problem, we get a dual problem:

$$\begin{aligned} \max_{\lambda, \alpha} \quad & \sum_{i=1}^n \lambda_i^2 \\ \text{subject to} \quad & \frac{1}{4\alpha^2} \left(\sum_{i=1}^n \lambda_i \mathbf{x}_i \right) \leq C \\ & \alpha \geq 0 \end{aligned}$$

This problem has an additional constraints than the original quadratic problem but different from ridge regression. It can use kernel trick as the solution contains inner product.