

Section 07: Kernels

1. Kernelized Linear Regression

Consider regularized linear regression (without a bias, for simplicity). We want to find the optimal parameters $\hat{w} = \arg \min_w L(w)$ for k -featured dataset $(x_i, y_i)_{i=1}^n$ (i.e. $X \in \mathbb{R}^{n \times k}$) that minimizes the following loss function:

$$L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

From class, we know there is an optimal \hat{w} that lies in the span of the datapoints. Concretely, there exists $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ such that $\hat{w} = \sum_i \alpha_i x_i$.

- Let $\alpha \in \mathbb{R}^n$, $\mathbf{K} \in \mathbb{R}^{n \times n}$, $\mathbf{y} \in \mathbb{R}^n$. Further, let us assume that we are using a linear kernel where $\mathbf{K}_{ij} = x_i^T x_j$. Express the loss function $L(w)$ in terms of \mathbf{K} and α as $L(\alpha)$.
- Assuming that \mathbf{K} is invertible, solve for the optimal $\hat{\alpha}$.
- Suppose after training our model on $\mathbf{X}_{\text{train}}$ we seek to make predictions on \mathbf{X}_{test} . Express these predictions $\hat{\mathbf{Y}}$ in terms of $\mathbf{K}_{\text{train}} = \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^T$, $\mathbf{y}_{\text{train}}$, $\mathbf{X}_{\text{train}}$, and \mathbf{X}_{test} . What would the general prediction formula look like if we are not using a linear kernel? Express the solution in terms of $\mathbf{K}_{\text{train, test}} = \Phi(\mathbf{X}_{\text{test}}) \Phi(\mathbf{X}_{\text{train}}^T)$ and $\hat{\alpha}$, where for arbitrary d , we have $\Phi: \mathbb{R}^k \rightarrow \mathbb{R}^d$.

2. Proving $\hat{w} \in \text{Span}(x_1, \dots, x_n)$

We will prove this through contradiction. Assume $\hat{w} \notin \text{Span}(x_1, \dots, x_n)$ solves $\arg \min_w L(w)$ where $L(w)$ is defined above. Then, there exists a component of \hat{w} that is perpendicular to the span, which we will call w^\perp . Concretely,

$$\hat{w} = \bar{w} + w^\perp$$

Where $\bar{w} = \sum_i \alpha_i x_i$ is the component of \hat{w} in the span of the datapoints.

- Show that $\hat{w} \cdot x_i = \bar{w} \cdot x_i$, for every x_i . (Hint: what is the relationship of w^\perp and x_i)
- Now, show that $\|\hat{w}\|_2^2 \geq \|\bar{w}\|_2^2$. (Hint: Use the Pythagorean Theorem)
- Finally, show that $\hat{w} \in \text{Span}(x_1, \dots, x_n)$. (Hint: Think about the regularization term)

3. Kernel Proofs

Let $\phi: d \rightarrow k$ be a feature map, and define K to be the kernel matrix of ϕ .

- Prove that the kernel matrix is symmetric. That is, show $K_{i,j} = K_{j,i}$.
- Recall that a matrix M is positive semi-definite if $x^T M x \geq 0, \forall x \in \mathbb{R}^n$. Show that K is positive semi-definite. (Hint: consider the matrix B where the i^{th} column of B is $\phi(x_i)$.)