

OLS Regression and PCA Result

Problem Setup: Let $X \in \mathcal{R}^{n \times d}$ be the features, $y \in \mathcal{R}^{n \times 1}$ be the labels. Assume X is full rank, implying that $\text{rank}(X) = \text{rank}(X^T X) = d$ and $n \geq d$.

Let the eigendecomposition of the covariance matrix be $X^T X = P^T \Lambda P$ i.e. the rows of $P \in \mathcal{R}^{d \times d}$ are the eigenvectors of the covariance matrix. Without loss of generality, assume that the rows are ordered with respect to the magnitude of the eigenvectors (the first row is the eigenvector corresponding to the largest eigenvalue, etc.) and that the eigenvectors are normalized, implying that $PP^T = I$.

Let $P_{k \times d}$ be the top- k rows of P with $k < d$. To perform Principal Component Analysis (PCA) on the data, compute $XP_{k \times d}^T$.

The ordinary-least squares (OLS) linear regression weights are given by:

$$w = (X^T X)^{-1} X^T y$$

Statement to Prove: We are concerned with Principal Component Regression (PCR) which involves reducing the dimensionality of a dataset using PCA, then learning an OLS model on the reduced data. We will show that the following two procedures result in the same learned weights:

1. (LHS) Learning weights in the original feature space, then applying PCA to the learned weights
2. (RHS) Learning OLS weights in the PCA-reduced feature space

Proof: We wish to show that $w_1 = w_2$ where

$$w_2 = P_{k \times d} (X^T X)^{-1} X^T y \quad \text{and} \quad w_1 = ((XP_{k \times d}^T)^T (XP_{k \times d}^T))^{-1} (XP_{k \times d}^T)^T y$$

Simplifying w_2 :

$$w_2 = P_{k \times d} (X^T X)^{-1} X^T y \tag{1}$$

$$= P_{k \times d} (P_{d \times d}^T \Lambda_{d \times d} P_{d \times d})^{-1} X^T y \tag{2}$$

$$= P_{k \times d} P_{d \times d}^T \Lambda_{d \times d}^{-1} P_{d \times d} X^T y \tag{3}$$

$$= I_{k \times d} \Lambda_{d \times d}^{-1} P_{d \times d} X^T y \tag{4}$$

$$= \Lambda_{k \times k}^{-1} P_{d \times d} X^T y \tag{5}$$

$$= \Lambda_{k \times k}^{-1} P_{k \times d} X^T y \tag{6}$$

Simplifying w_1 :

$$w_1 = ((XP_{k \times d}^T)^T (XP_{k \times d}^T))^{-1} (XP_{k \times d}^T)^T y \quad (7)$$

$$= (P_{k \times d} X^T X P_{k \times d}^T)^{-1} P_{k \times d} X^T y \quad (8)$$

$$= (P_{k \times d} P_{d \times d}^T \Lambda_{d \times d} P_{d \times d} P_{k \times d}^T)^{-1} P_{k \times d} X^T y \quad (9)$$

$$= (I_{k \times d} \Lambda_{d \times d} I_{d \times k})^{-1} P_{k \times d} X^T y \quad (10)$$

$$= \Lambda_{k \times k}^{-1} P_{k \times d} X^T y \quad (11)$$

Therefore $w_1 = w_2 = \Lambda_{k \times k}^{-1} P_{k \times d} X^T y$, and both procedures result in the same learned weights.