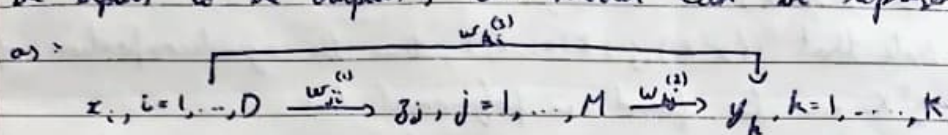


## SET-6

6.1) Deep Learning2) Exercise 5.18 (Bishop (2006))

- \* On adding a new ship layer with connections directly from the inputs to the outputs, our model can be represented as:



- \* The equations in Section 5.3.2 mostly remain the same, except Eq. (5.64), which is now

$$y_k = \sum_{j=1}^M w_{jk}^{(2)} z_j + \sum_{i=1}^D w_{ki}^{(3)} x_i$$

- \* The eq. for the derivative of the error function w.r.t.  $w_{ki}^{(3)}$  is

$$\begin{aligned} \frac{\partial E_n}{\partial w_{ki}^{(3)}} &= (y_n - t_n) \frac{\partial y_k}{\partial w_{ki}^{(3)}} = (y_n - t_n) x_i \\ &\equiv \delta_k x_i \end{aligned}$$

6.2) Representation Learning1) a) Exercise 23.1

- 1) \* We know that if  $V, W$  are finite dimensional vector spaces, &  $T$  is a linear transformation from  $V$  to  $W$ , then the image of  $T$  is a finite-dimensional subspace of  $W$  and

$$\dim(V) = \dim(\text{null}(T)) + \dim(\text{image}(T))$$

- \*  $\therefore$  for  $\dim(\text{null}(A)) \geq 1$ ,  $\exists v \neq 0$  s.t.  $Av = A_0 = 0$ .

$$\therefore \exists u \neq v \in \mathbb{R}^n \text{ s.t. } Au = Av$$

- 2) We have that  $Au = Av$  for some  $u \neq v \in \mathbb{R}^n$ .

$$\Rightarrow \text{For any recovery function } f, f(Au) = f(Av)$$

- $\therefore$  Exact recovery of a linear compression scheme is impossible.

b) Exercise 23-3

- \* Let  $X$  be a matrix with its  $j$ -th column as  $\psi(x_j)$ .
- \* We find the spectral decomposition of  $X^T X$ , using the results from section 23.1.1 (A more efficient solution for the case  $d \gg m$ ).
- \* Note that  $(X^T X)_{ij} = K(x_i, x_j)$ , thus the eigendecomposition of  $X^T X$  can be found in polynomial time.
- \* Let  $V$  be the matrix whose columns are the  $n$  leading eigenvectors of  $X^T X$ , and let  $D$  be a diagonal  $n \times n$  matrix whose diagonal consists of the corresponding eigenvalues.
- \* Denote by  $U$  be the matrix whose columns are the  $n$  leading eigenvectors of  $XX^T$ . ~~We now show how to~~
- \* We now show how to project the data without maintaining the matrix  $U$ . For every  $x \in X$ , the projection  $U^T \phi(x)$  is calculated as follows:

$$U^T \phi(x) = D^{\frac{1}{2}} V^T X^T \phi(x) = D^{\frac{1}{2}} V^T \begin{pmatrix} K(x_1, x) \\ \vdots \\ K(x_m, x) \end{pmatrix}$$



c) Exercise 23.4

1) \* We have that  $\forall w \in \mathbb{R}^d$  s.t.  $\|w\|=1$  & every  $i \in [m]$ ,  
 $(\langle w, x_i \rangle)^2 = \text{tr}(w^T x_i x_i^T w)$

\* Hence

$$\begin{aligned} \underset{w: \|w\|=1}{\text{argmax}} \text{Var}[\langle w, x \rangle] &= \underset{w: \|w\|=1}{\text{argmax}} \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle)^2 \\ &= \underset{w: \|w\|=1}{\text{argmax}} \frac{1}{m} \sum_{i=1}^m \text{tr}(w^T x_i x_i^T w) \end{aligned}$$

\* We see that our equation reduces to the PCA problem for  $n=1$ .

\* Hence the optimal solution of our the variance maximization problem is to set  $w$  to be the first principle vector of  $x_1, \dots, x_m$ .

2) \* We know that

$$E[(\langle w, x \rangle)(\langle w, x \rangle)] = w^T E[x x^T] w = m w^T A w, \text{ where } A = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

Since  $w$  is an eigen vector of  $A$  and  $E[(\langle w, x \rangle)(\langle w, x \rangle)] = 0$

$$\Leftrightarrow \langle w, w \rangle = 0$$

\*  $\therefore$  Our problem reduces to

$$w^* = \underset{w: \|w\|=1, E[(\langle w, x \rangle)(\langle w, x \rangle)] = 0}{\text{argmax}} \text{Var}[\langle w, x \rangle]$$

$$= \underset{w: \|w\|=1, \langle w, w \rangle = 0}{\text{argmax}} \text{Var}[\langle w, x \rangle]$$

$$\Rightarrow w^* = \underset{w: \|w\|=1, \langle w, w \rangle = 0}{\text{argmax}} \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle)^2$$

$$= \underset{w: \|w\|=1, \langle w, w \rangle = 0}{\text{argmax}} \text{tr}(w^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w)$$

\* Note that the PCA problem for  $n=2$  is equivalent to finding a unitary matrix  $W \in \mathbb{R}^{d \times 2}$  such that

$$(W^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T W) \text{ is maximized.}$$

\* We know that the columns of the optimal matrix,  $w_1, w_2$  are the two first principal vectors of  $x_1, \dots, x_m$ .

$$\Rightarrow W^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T W = w_1^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w_1 + w_2^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w_2$$

\* Since  $w^*$  and  $w_1$  are orthonormal, we obtain that

$$w_1^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w_1 + w_2^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w_2 \geq w_1^T \frac{1}{m} \sum_{i=1}^m x_i x_i^T w_1 + w^{*T} \frac{1}{m} \sum_{i=1}^m x_i x_i^T w^*$$

$\therefore$  we conclude that  $w^* = w_2$ .

20.5) Exercise 20.5

a) Covariance of the deflated matrix is given by

$$\tilde{C} = \frac{1}{n} ((I - v_i v_i^T) X^T X (I - v_i v_i^T))$$

$$= \frac{1}{n} (X^T X - v_i v_i^T X^T X) (I - v_i v_i^T)$$

$$= \frac{1}{n} ((X^T X - v_i n \lambda_i v_i^T) (I - v_i v_i^T))$$

$$= \frac{1}{n} (X^T X - X^T X v_i v_i^T - v_i n \lambda_i v_i^T + v_i n \lambda_i v_i^T v_i v_i^T)$$

$$= \frac{1}{n} (X^T X - n \lambda_i v_i v_i^T)$$

$$= \frac{1}{n} X^T X - \lambda_i v_i v_i^T$$

b) As  $\tilde{X} \in (d-1)$  subspace orthogonal to  $v_i$ ,

$\Rightarrow u$  must be orthogonal to  $v_i$

$$\therefore \Rightarrow u^T v_i = 0 \quad \& \quad u^T u = 1$$

$$\therefore u = v_i$$

c) def function  $[V, \lambda] = \text{simple PCA}(C, K, f)$  &

$$d = C.\text{length}();$$

$$V = \text{zeros}(d, K);$$

for  $j$  in range  $(1, K)$  &

$$[\lambda(j), v(:,j)] = f(C);$$

$$C = C - \lambda(j) * v(:,j) * v(:,j)^T;$$

}

}