

## Ex 2.1

- a) A linear dataset can be simplified with a simple linear equation where it has a polynomial degree of 1. This means two variables can be linearly separated with a line unlike non-linear data which has many dimensions and is difficult to visualize.

PCA can easily comprehend this linear relationship in the data however if the data has nonlinear pattern, PCA might reverse the analysis and point to the exact opposite direction. This can also result to information loss along with biased output. Therefore, it is important to make sure the data is correct.

- b) PCA essentially maximizes the variance to identify which components contribute how much in our dataset. If we do not normalize and data points are very far away from each other it is likely that PCA will take the largest variance as the most important component and provide a biased output. For example, in this case we might not be able to see the contribution of data points at all.

- but one of other can be highlighted.  
Only one feature may get highlighted.  
This is why normalizing the data is important where all features are taken into account.

c) i) PCA does not take account of label information. Therefore, always reducing dimensions may lead to losing important features. For example, when distinguishing between cars and buses, taking how many miles they have gone factor won't work even though it may have large variance. Height might be a good factor to consider. If the dataset is arranged in a way that PCA is biased towards wheels discarding height as a factor would be wrong.

ii) When the data is in different scales or notations, PCA should not be applied. Such as if we have categorical data or numerical information that are in different scales (i.e. currency), PCA would result in wrong components and features.

we would lose important information by reducing dimensions.

iii) Another case would be health related data where data might not fluctuate as much thus resulting in low variance but could be an important factor. Blindly reducing dimension here would cause information loss and ineffective classification.

#### Ex. 2.4

a) Given,

$$y = x^T \cdot w + \epsilon$$

Assuming  $\epsilon$  as normally distributed,

$$f(y) = x^T \cdot w$$

Conditional mean of  $y$  given  $x$ ,

$$P(y|x) = x^T \cdot w$$

and the conditional variance of  $y$  given  $x$  is

$$\text{Var}(y) = \sigma^2$$

... probability function of  $y$  given  $x$ ,

Conditional probability:

$$P(y|x) = \frac{P(\{x^T \cdot w\} \cap y)}{P(x^T \cdot w)}$$

b) Since  $\{x_n, y_n\}_{n=1}^N$ ;  $x_n \in \mathbb{R}^D$   
hence the mode becomes:

$$y_i = x_i^T \cdot w + \epsilon \quad ; \text{ where } i=1, 2, \dots, n$$

So the likelihood function is:

$$L(x_i, y_i; w, \sigma^2, \epsilon) = \prod_{i=1}^n \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left[ -\frac{1}{2\sigma^2} (y_i - x_i^T w - \epsilon)^2 \right]$$

d) To calculate MLE for  $w$ :

Removing the constants and taking  
log of the function

$$L = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i^T w - \epsilon)^2$$

To find the MLE's for  $w$ , we have to find LSE of  $w$ .

$$\frac{\partial L}{\partial \epsilon} = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2 (y_i - x_i^T w - \epsilon)(-1) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - n\epsilon - w \sum_{i=1}^n x_i^T = 0$$

$$\Rightarrow \epsilon = \bar{y} - w \bar{x}^T$$

$$\frac{\partial L}{\partial w} = \frac{1}{-2\sigma^2} \sum_{i=1}^n 2 (y_i - x_i^T w - \epsilon)(-x_i)$$

$$\Rightarrow \sum_{i=1}^n y_i x_i^T - \epsilon \sum_{i=1}^n x_i^T - w \sum_{i=1}^n x_i^T x_i^T = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i^T - (\bar{y} - w \bar{x}^T) \sum_{i=1}^n x_i^T - w \sum_{i=1}^n x_i^T x_i^T = 0$$

$$\Rightarrow w = \frac{\sum_{i=1}^n (x_i^T - \bar{x}^T)(y_i - \bar{y})}{\sum_{i=1}^n (x_i^T - \bar{x}^T)^2}$$

This is the maximum likelihood estimator for  $w$ .

