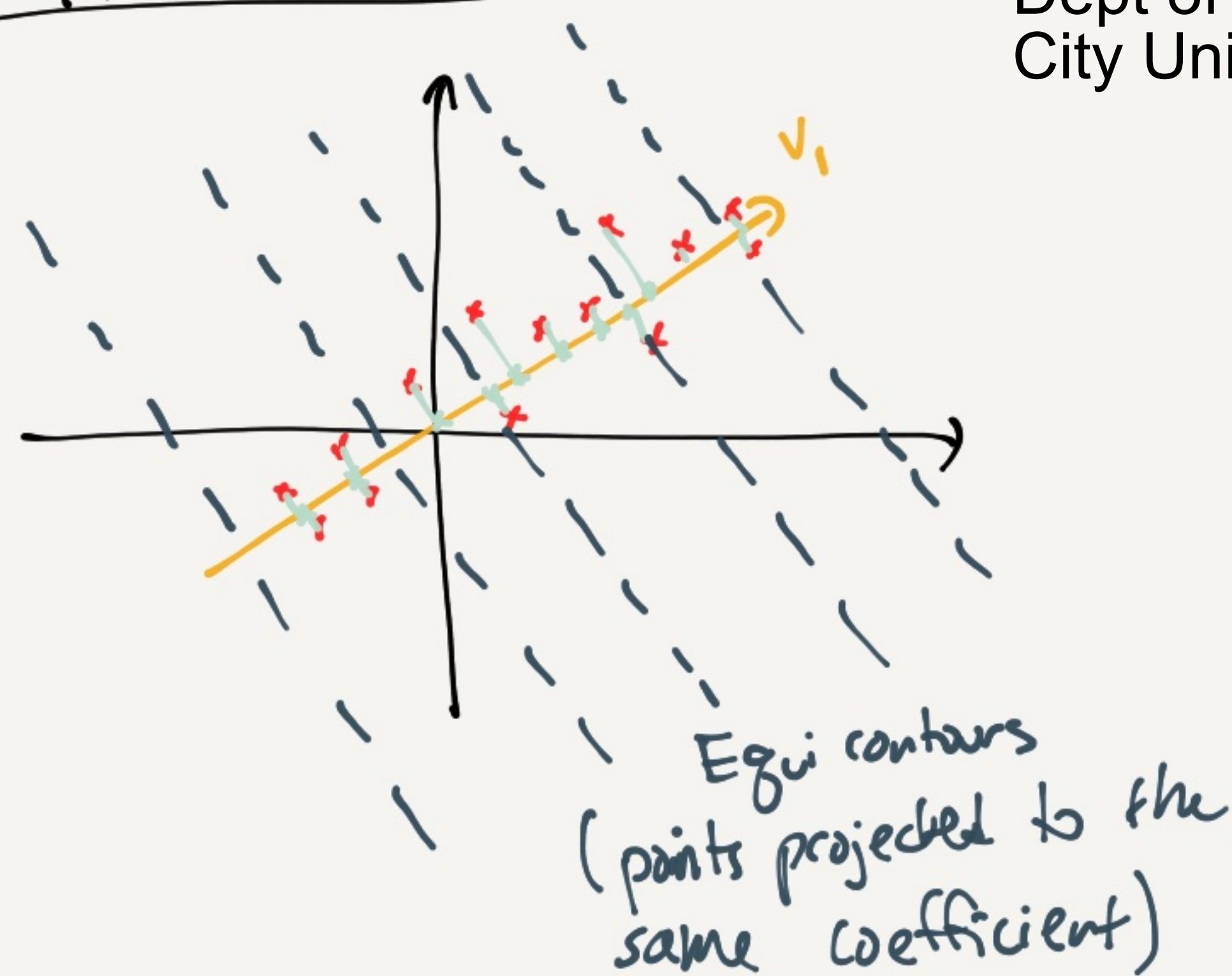
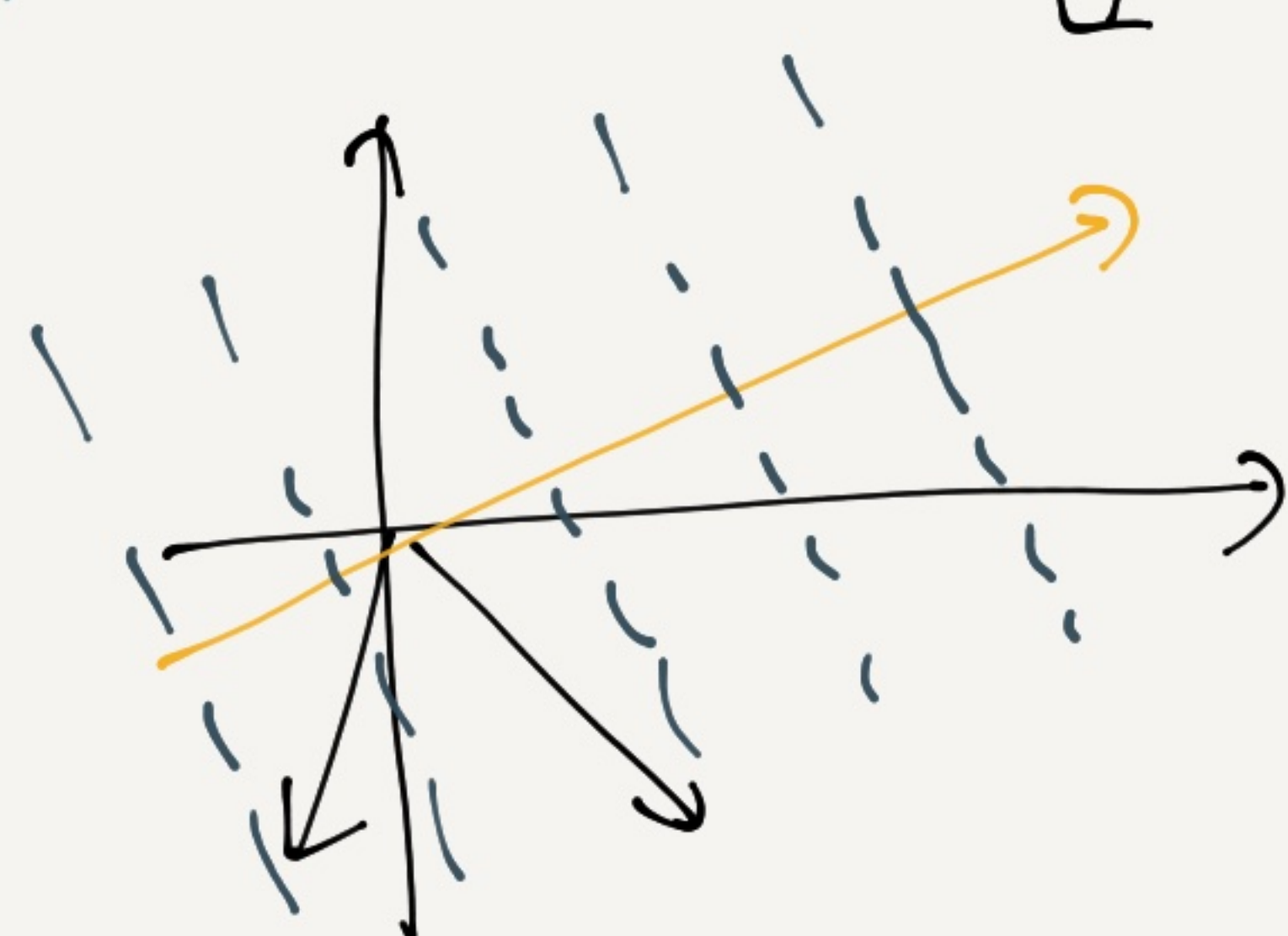
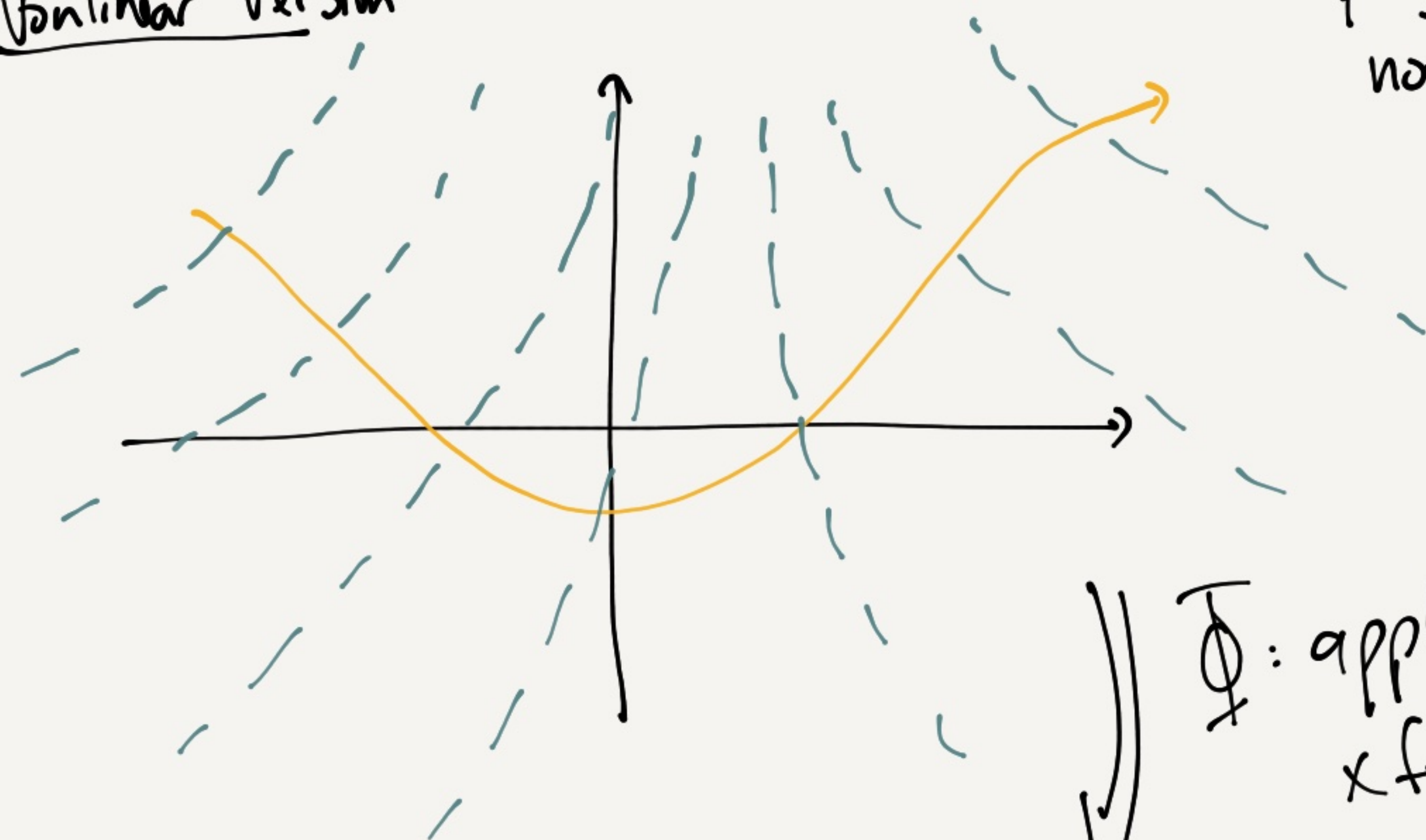


# Lecture 11 - Non linear dimensionality Reduction

## PCA (linear dim reduction)



## Nonlinear Version



CS5487 Lecture Notes (2020)  
Dr. Antoni B. Chan  
Dept of Computer Science  
City University of Hong Kong

project onto  
nonlinear function.  
Assignment Project Exam Help  
<https://powcoder.com>  
Add WeChat powcoder

$\Phi$ : apply feature  
x form

linear xform  
in high-dim  
space.

## Kernel PCA

• Data:  $\{x_1, \dots, x_n\}$

• Apply feature  $x_{form}$  to  $x_i \rightarrow \phi(x_i)$   
 $X \rightarrow \Phi$

• Assume the  $x_{form}$  data is zero-mean (centered)  
i.e.  $\sum_i \phi(x_i) = 0$

• Covariance matrix: (assuming zero-mean)  $C = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T$

• Find Eigenvectors/values:  $C v_j = \lambda_j v_j$   
( $v_j, \lambda_j$ )

$$\frac{1}{n} \sum_{m=1}^n \phi(x_m) \phi(x_m)^T v_j = \lambda_j v_j$$

scalar  $a_{j,m}$

$$\lambda_j v_j = \frac{1}{n} \sum_{m=1}^n a_{j,m} \phi(x_m)$$

Thus  $v_j$  has the form:

$$v_j = \sum_i a_{j,i} \phi(x_i) \quad \text{for some coefficients } a_{j,i}$$

$$= \Phi a_j, \quad a_j = \begin{bmatrix} a_{j,1} \\ \vdots \\ a_{j,n} \end{bmatrix}$$

Substitute for  $v_j$

$$\Rightarrow \frac{1}{n} \sum_m \phi(x_m) \phi(x_m)^T \Phi a_j = \lambda_j \Phi a_j$$

$\Phi \Phi^T$

$$\Rightarrow \frac{1}{n} \Phi \Phi^T \Phi a_j = \lambda_j \Phi a_j$$



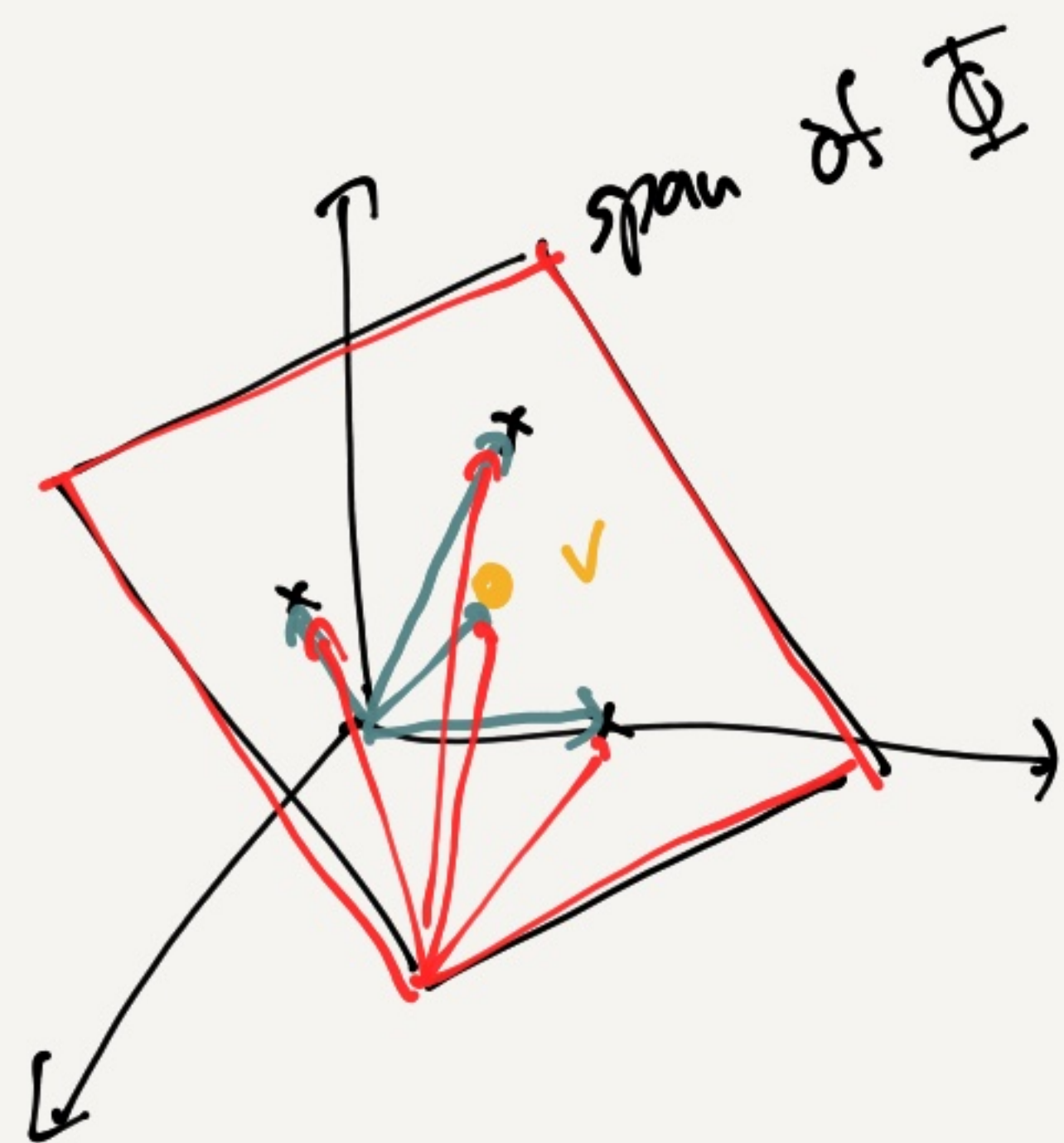
$$\frac{1}{n} \Phi \Phi^T \Phi a_j = \lambda_j \Phi a_j$$

eigenvector is a linear combo of  $\phi(x_1) \dots \phi(x_n)$

Pre-multiply by  $\Phi^T$ :

(equivalent set of equations by projecting onto the span of  $\Phi$ )

$$\Rightarrow \frac{1}{n} \Phi^T \Phi \Phi^T \Phi a_j = \lambda_j \Phi^T \Phi a_j$$



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Apply kernel trick

$$\frac{1}{n} K K a_j = \lambda_j K a_j$$

$$\Rightarrow K K a_j = n \lambda_j K a_j \quad (1)$$

Two possibilities:

1)  $K$  is invertible  $\Rightarrow$  multiply by  $K^{-1}$

$$\Rightarrow \text{solve } K a_j = n \lambda_j a_j$$

2)  $K$  is not invertible  $\Rightarrow$  the only difference is eigenvectors with  $\lambda_j = 0 \Rightarrow$  do not affect the PCA solution.

$$\Rightarrow \text{solve } K a_j = n \lambda_j a_j$$

Verification

Suppose  $(a_j, \lambda_j)$  are an eigenpair of  $K$  s.t.

$$\checkmark K a_j = n \lambda_j a_j, \quad \lambda_j \neq 0, \quad a_j^T a_j = 1$$

Then looking at (1):

$$K K a_j = n \lambda_j \underbrace{K a_j}_{n \lambda_j a_j} \quad \checkmark$$

$$\underbrace{n \lambda_j K a_j}_{n \lambda_j a_j} = n^2 \lambda_j^2 a_j$$

$$n^2 \lambda_j^2 a_j = n^2 \lambda_j^2 a_j \quad \checkmark$$



PC  $v_j$  should be normalized

$$v_j^T v_j = 1$$

$$a_j^T \Phi^T \Phi a_j = 1$$

$$a_j^T \underbrace{K}_{n \times n} a_j = 1$$

$$\underbrace{n \lambda_j}_{\text{eigenvalue}} a_j^T a_j = 1$$

$$\Rightarrow a_j \leftarrow \frac{1}{\sqrt{n \lambda_j}} a_j$$

$$\text{Then } v_j^T v_j = 1$$

Kernel Centering - How to center the data in the feature space  $\phi(x)$ ?

Consider  $x_i$ , to center  $\phi(x_i)$ :

$$\tilde{\phi}(x_i) = \phi(x_i) - \frac{1}{n} \sum_j \phi(x_j) = \phi(x_i) - \frac{1}{n} \Phi^T \mathbf{1}$$

Center all data:

$$\tilde{\Phi} = \Phi - \left[ \frac{1}{n} \Phi^T \mathbf{1} \right] \mathbf{1}^T = \Phi \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)$$

Centered Kernel:

$$\tilde{K} = \tilde{\Phi}^T \tilde{\Phi} = \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)^T \Phi^T \Phi \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)$$
$$\boxed{\tilde{K} = \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)^T K \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)}$$

## Summary KPCA (Kernel PCA)

Given data  $X$ :

- training*
- 1) calculate kernel matrix:  $K = [k(x_i, x_j)]_{i,j}$
  - 2) center kernel:  $\tilde{K} = \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) K \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)$
  - 3) Find the top  $D$  eigenvectors of  $\tilde{K}$ :  $\tilde{K} a_j = \lambda_j a_j$
  - 4) scale  $a_j \leftarrow \frac{1}{\sqrt{\lambda_j}} a_j$

- testing*
- 5) project  $x_*$ :  $z_{*,j} = \tilde{\phi}(x_*)^T \tilde{\Phi} a_j = \tilde{K}_*^T a_j$
- $$\tilde{K}_* = K_* - \frac{1}{n} K_*^T \mathbf{1} - \frac{1}{n} K \mathbf{1} + \frac{1}{n^2} \mathbf{1}^T K \mathbf{1}$$

$$K_* = [k(x_*, x_i)]_i \leftarrow \text{kernel function between new point } x_* \text{ \& data } X.$$

• Same properties of PCA in  $\mathcal{H}$

- The first  $D$  PC carry the most variance (in  $\mathcal{H}$ )  
minimize the reconstruction error (in  $\mathcal{H}$ )
- PC are orthogonal (in  $\mathcal{H}$ )

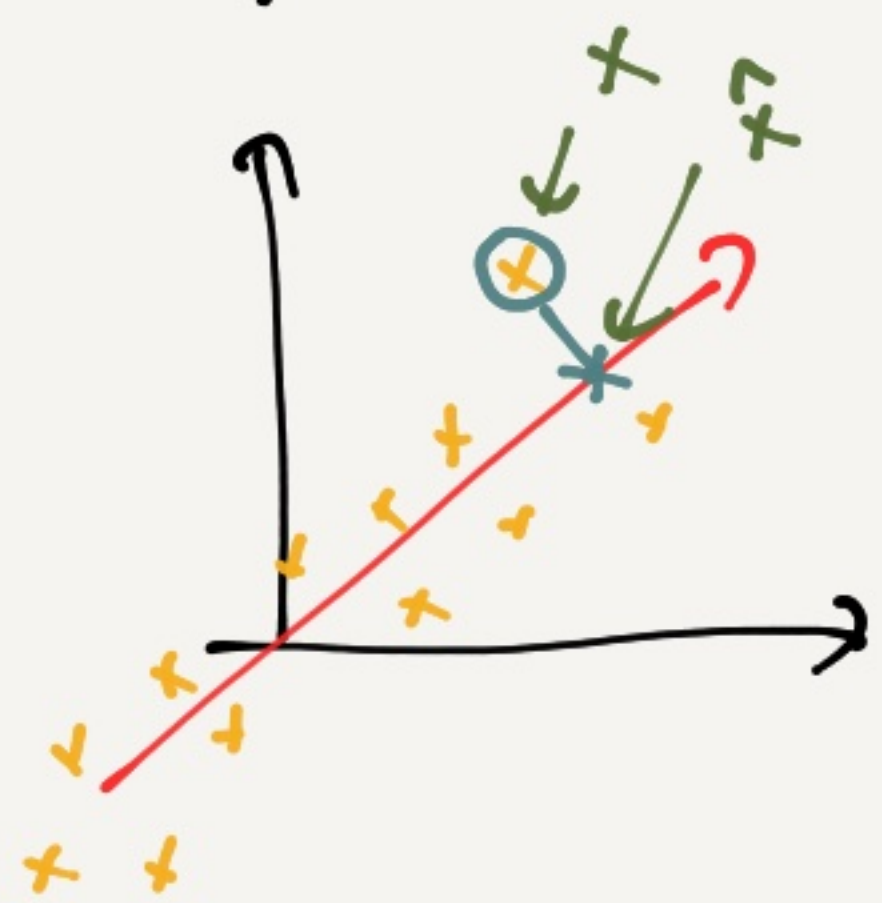
Final: the original problem needs a  $d$ -dim eigenvector  
 $\Rightarrow$  kernelized problem needs a  $N$ -dim eigenvector.

• Kernel version of FLD  $\Rightarrow$  Kernel discriminant analysis (10.15)



## Pre-image Problem

Given PCA coefficients  $z$ , we can reconstruct  $\hat{x}$ .  
(e.g. denoising)



$$\text{given } z, \quad \hat{x} = \sum_j v_j z_j = Vz$$

## What about KPCA?

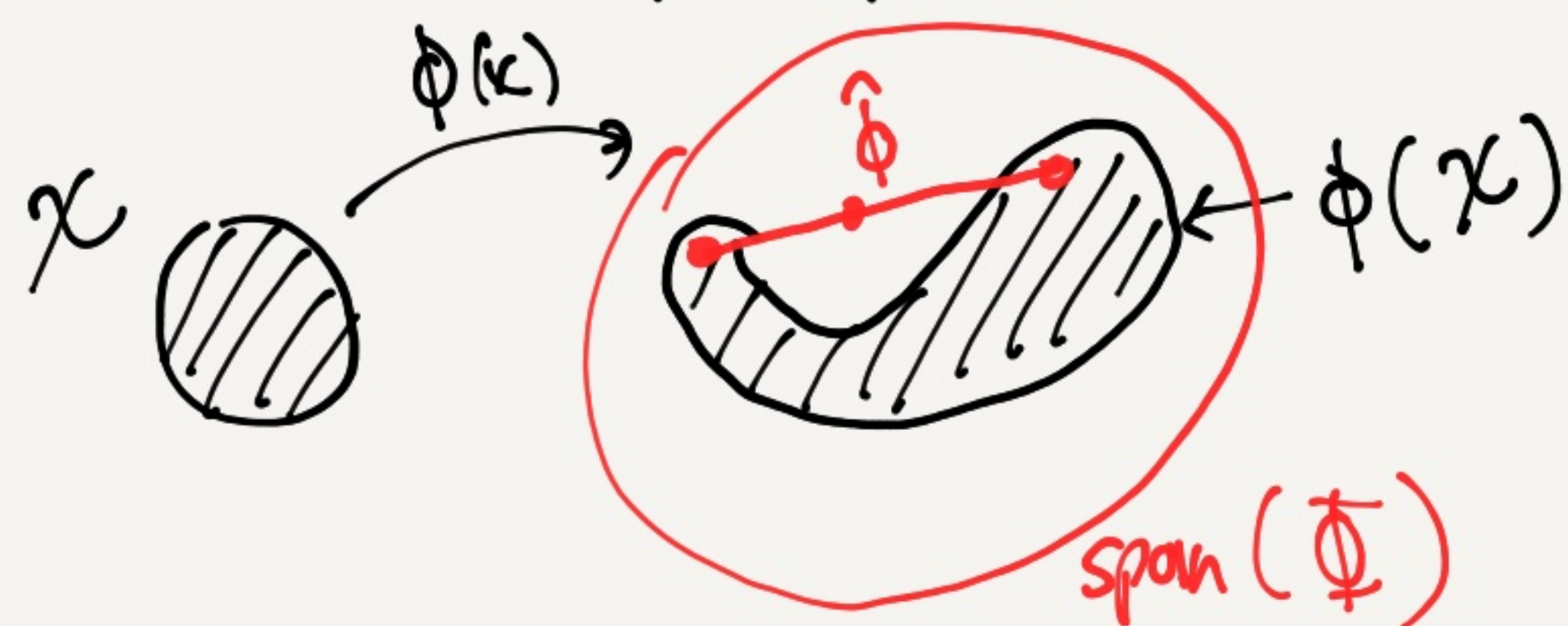
• given a  $z$ , what is the high-dim feature?

$$\hat{\phi} = \sum_j v_j z_j = \sum_j \Phi a_j z_j = \Phi \left( \sum_j a_j z_j \right)$$

linear combo of  $\phi(x_i)$ 's  
 $\hat{\phi} \in \text{span}(\Phi)$

← reconstruction

• But... not all points in  $\text{span}(\Phi)$  have a corresponding  $x$  in the input space.



Example: let  $\phi(x) = \begin{bmatrix} x^2 \\ x \end{bmatrix}$

$$x_1 = 1 \Rightarrow \phi(x_1) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$x_2 = 2 \Rightarrow \phi(x_2) = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$

$$\hat{\phi} = \phi(x_1) + \phi(x_2) = \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

What is the  $\hat{x}$  that gives  $\hat{\phi}$ ?  
There's no  $\phi(\hat{x}) = \hat{\phi}$

## Approximate Pre-image

• Consider the general problem:  $\hat{\phi} = \Phi \alpha$  (given vector)

• Goal: find  $\hat{x}$  that gives the closest  $\phi(\hat{x})$  to  $\hat{\phi}$ .  
minimum norm reconstruction.

$$\hat{x} = \underset{x \in X}{\operatorname{argmin}} \|\phi(x) - \hat{\phi}\|^2$$

$$= \underset{x \in X}{\operatorname{argmin}} K(x, x) - 2 \sum_i \alpha_i k(x, x_i) + \underbrace{\sum_i \sum_j \alpha_i \alpha_j k(x_i, x_j)}_{\text{constant wrt } x}$$

$$\hat{x} = \underset{x \in X}{\operatorname{argmin}} K(x, x) - 2 \sum_i \alpha_i k(x, x_i)$$

Soln 1: (nearest neighbor): select  $x \in X$  (dataset)

$$\hat{x} = \underset{j}{\operatorname{argmin}} k(x_j, x_j) - 2 k_j^T \alpha$$

where  $k_j = j^{\text{th}}$  column of kernel matrix  $K$

Soln 2: solve the minimization directly w/ optimizer.

Soln 3: Suppose  $k(x, x) = 1$ ,  $k(x_i, x_j) \geq 0$  (e.g. Gaussian kernel) and  $k$  is homogeneous,  $k(\|x_i - x_j\|)$

$$\hat{x} = \underset{x}{\operatorname{argmin}} - \sum_i \alpha_i k(x_i, x) = \underset{x}{\operatorname{argmax}} \sum_i \alpha_i k(x_i, x) = \underset{x}{\operatorname{argmax}} \left( \sum_i \alpha_i k(x_i, x) \right)^2$$

$\Rightarrow$  iterative algorithm:  
(like mean-shift w/ weights  $\alpha_i$  on each  $x_i$ )

$$\hat{x} \leftarrow \frac{\sum_i \alpha_i k'(\|x_i - \hat{x}\|^2) x_i}{\sum_i \alpha_i k'(\|x_i - \hat{x}\|^2)}$$