CS2109S Tutorial 7 Unsupervised Learning

(AY 24/25 Semester 2)

March 27, 2025

(Prepared by Benson)

## Contents

#### Clustering

- Q1. K-Means Algorithm
- Q2. Kernel K-Means

### Principal Component Analysis (PCA)

Q3. Lossy Compression Bonus. Composing Kernels

#### Intuition:

- Lines 4-5: Minimize the distortion by reassigning clusters, keeping the centroids unchanged.
- Lines 6-7: Minimize the distortion by changing centroids, keeping the assignment unchanged.

#### Algorithm K-Means Clustering

1: for k = 1 to K do 2:  $\mu_k \leftarrow \text{random location}$ 3: while not converged do 4: for i = 1 to m do5:  $c^{(i)} \leftarrow \operatorname{argmin}_k \|\mathbf{x}^{(i)} - \mu_k\|^2$ 6: for k = 1 to K do7:  $\mu_k \leftarrow \frac{1}{|\{\mathbf{x}^{(i)}|c^{(i)}=k\}|} \sum_{\mathbf{x} \in \{\mathbf{x}^{(i)}|c^{(i)}=k\}} \mathbf{x}$ 

$$\frac{d\mathcal{E}}{d\boldsymbol{\mu}_k} = 0 \quad \Rightarrow \quad \sum_{\boldsymbol{c}^{(i)}=k} 2(\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k) = 0 \quad \Rightarrow \quad \boldsymbol{\mu}_k = \frac{1}{|\{\boldsymbol{x}^{(i)}|\boldsymbol{c}^{(i)}=k\}|} \sum_{\boldsymbol{c}^{(i)}=k} \boldsymbol{x}^{(i)}$$

- (a) Given that every iteration produces a partition with a lower distortion, prove that this algorithm always converges. Convergence is when the centroids/medoids do not change after an iteration of the algorithm.
- There are not more than k<sup>N</sup> possible partitions.
- Since every iteration of the algorithm produces a partition that has a lower distortion, it eventually terminates.

Algorithm K-Means Clustering					
1: for $k = 1$ to $K$ do					
2: $\mu_k \leftarrow$ random location					
3: while not converged do					
4: <b>for</b> $i = 1$ <b>to</b> $m$ <b>do</b>					
5: $   c^{(i)} \leftarrow \operatorname{argmin}_k \  oldsymbol{x}^{(i)} - oldsymbol{\mu}_k \ ^2 $					
6: for $k = 1$ to $K$ do					
7: $\left[ \begin{array}{c} \mu_k \leftarrow rac{1}{ \{m{x}^{(i)} m{c}^{(i)}=k\} } \sum_{m{x}\in\{m{x}^{(i)} m{c}^{(i)}=k\}} m{x} \end{array} \right]$					

(b) Although k-means always converges, it may get stuck at a bad local minimum. As mentioned in the lecture, one method of circumventing this is to run the algorithm multiple times and choose the clusters with the minimum distortion. Suggest some other ways to help the algorithm get closer to the global minimum.



Intuition: Pick centroids that are as far as possible during random initialization.

Determinstic method:

Algorithm Cluster initialization

- 1:  $\mu_1 \leftarrow$  random location
- 2: for k = 2 to K do
- 3:  $\mu_k \leftarrow \text{ point farthest possible from } \mu_1, \dots, \mu_{k-1}$

Random method (K-means++):

#### Algorithm Cluster initialization

- 1:  $\mu_1 \leftarrow$  random location
- 2: for k = 2 to K do
- 3:  $D(x) \leftarrow \text{distance from } x \text{ to closest centroid in } \mu_1, \dots, \mu_{k-1}$
- 4:  $\mu_k \leftarrow$  random point with (weighted) probability distribution  $D(x)^2$

#### K-means++ Initialization



 (c) Cluster the 6 points in table 1 into two clusters using the K-means algorithm. The two initial centroids are (0,1) and (2.5,2).

#### Iteration 1

i	1	2	3	4	5	6
x	1	1	2	2	3	3
y	0	1	1	2	1	2
С	1	1	2	2	2	2

• 
$$\mu_1 = \frac{1}{2}((1,0) + (1,1)) = (1,0.5)$$
  
•  $\mu_2 = \frac{1}{4}((2,1) + (2,2) + (3,1) + (3,2)) = (2.5,1.5)$ 

Algorithm K-Means Clustering 1: for k = 1 to K do 2:  $\mu_k \leftarrow$  random location 3: while not converged do for i = 1 to m do 4.  $c^{(i)} \leftarrow \operatorname{argmin}_{k} \| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{k} \|^{2}$ 5: for k = 1 to K do 6·  $\mu_k \leftarrow \frac{1}{|\{\mathbf{x}^{(i)}| c^{(i)} = k\}|} \sum_{\mathbf{x} \in \{\mathbf{x}^{(i)}| c^{(i)} = k\}} \mathbf{x}$ 7: 2  $\succ$ 0 0 2 3 x

 (c) Cluster the 6 points in table 1 into two clusters using the K-means algorithm. The two initial centroids are (0,1) and (2.5,2).

#### Iteration 2

i	1	2	3	4	5	6
x	1	1	2	2	3	3
y	0	1	1	2	1	2
С	1	1	2	2	2	2

• 
$$\mu_1 = \frac{1}{2}((1,0) + (1,1)) = (1,0.5)$$
  
•  $\mu_2 = \frac{1}{4}((2,1) + (2,2) + (3,1) + (3,2)) = (2.5,1.5)$ 

Algorithm K-Means Clustering 1: for k = 1 to K do 2:  $\mu_k \leftarrow$  random location 3: while not converged do for i = 1 to m do 4.  $| c^{(i)} \leftarrow \operatorname{argmin}_{k} \| \mathbf{x}^{(i)} - \boldsymbol{\mu}_{k} \|^{2}$ 5: for k = 1 to K do 6·  $\mu_k \leftarrow rac{1}{|\{m{x}^{(i)}| c^{(i)}=k\}|} \sum_{m{x}\in\{m{x}^{(i)}| c^{(i)}=k\}} m{x}$ 7: 2 ×  $\succ$ × 0 0 2 3 x

(d) Cluster the 6 points in table 1 into two clusters using the K-medoids algorithm. The initial medoids are point 1 and point 3.

#### Iteration 1

i	1	2	3	4	5	6
x	1	1	2	2	3	3
У	0	1	1	2	1	2
С	1	2*	2	2	2	2

• 
$$\mu_1 = (1,0)$$
  
•  $\mu_2 = \frac{1}{5}((1,1) + (2,1) + (2,2) + (3,1) + (3,2)) = (2.2,1.4) \Rightarrow (2,1)$ 

#### Algorithm K-Medoids Clustering

1. for k = 1 to K do 2:  $\mu_k \leftarrow \text{random data point } \mathbf{x}^{(i)}$ 3: while not converged do for i = 1 to m do Δ·  $c^{(i)} \leftarrow \operatorname{argmin}_{k} \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{k} \|^{2}$ 5. for k = 1 to K do 6:  $oldsymbol{\mu}_k \leftarrow rac{1}{|\{oldsymbol{x}^{(i)}| c^{(i)}=k\}|} \sum_{oldsymbol{x} \in \{oldsymbol{x}^{(i)}| c^{(i)}=k\}} oldsymbol{x}$ 7: for k = 1 to K do  $\triangleright$  Closest data point 8:  $\boldsymbol{\mu}_k \leftarrow \operatorname{argmin}_{\boldsymbol{x} \in \{\boldsymbol{x}^{(i)} | \boldsymbol{c}^{(i)} = k\}} \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k \|^2$ 9: 2 Closest data point  $\succ$ 0 2 3 n х

10 / 32

## Recap: Kernel Trick

Important Property of a Kernel Function  $K(\boldsymbol{u}, \boldsymbol{v}) = \phi(\boldsymbol{u}) \cdot \phi(\boldsymbol{v}) \text{ for some function } \phi$ 

If we have a model that **ONLY** relies on the dot products of  $x^{(i)}$ :

$$\mathsf{SVM}_{oldsymbol{lpha}}(oldsymbol{x}) = \mathsf{sign}\left(\sum_{i=1}^N lpha^{(i)} y^{(i)}(oldsymbol{x}^{(i)} \cdot oldsymbol{x})
ight)$$

Applying the kernel function  $\equiv$  applying the transformed features!

$$\mathsf{SVM}_{\alpha}(\mathbf{x}) = \mathsf{sign}\left(\sum_{i=1}^{N} \alpha^{(i)} y^{(i)}(\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}))\right) = \mathsf{sign}\left(\sum_{i=1}^{N} \alpha^{(i)} y^{(i)} \mathcal{K}(\mathbf{x}^{(i)}, \mathbf{x})\right)$$

(a) Clustering 1 should be considered better than Clustering 2. Can the K-means algorithm achieve the better clustering? Why or why not?



- (a) Clustering 1 should be considered better than Clustering 2. Can the K-means algorithm achieve the better clustering? Why or why not?
  - No. K-means uses **Euclidean distance**.
  - Assignment to clusters will follow a linear boundary.



(b) Write the squared Euclidean distance between two points  $p_i$  and  $p_j$  using vector norms (length) and the dot product. How can we apply the kernel trick?

$$\begin{split} \|\boldsymbol{p}_i - \boldsymbol{p}_j\|^2 &= (\boldsymbol{p}_i - \boldsymbol{p}_j) \cdot (\boldsymbol{p}_i - \boldsymbol{p}_j) \\ &= \boldsymbol{p}_i \cdot \boldsymbol{p}_i - \boldsymbol{p}_i \cdot \boldsymbol{p}_j - \boldsymbol{p}_j \cdot \boldsymbol{p}_i + \boldsymbol{p}_j \cdot \boldsymbol{p}_j \\ &= \boldsymbol{p}_i \cdot \boldsymbol{p}_i - 2(\boldsymbol{p}_i \cdot \boldsymbol{p}_j) + \boldsymbol{p}_j \cdot \boldsymbol{p}_j \end{split}$$

Naively applying transformed features:

$$\phi(\boldsymbol{p}_i) \cdot \phi(\boldsymbol{p}_i) - 2(\phi(\boldsymbol{p}_i) \cdot \phi(\boldsymbol{p}_j)) + \phi(\boldsymbol{p}_j) \cdot \phi(\boldsymbol{p}_j)$$

Kernel trick:

$$K(\boldsymbol{p}_i, \boldsymbol{p}_i) - 2K(\boldsymbol{p}_i, \boldsymbol{p}_j) + K(\boldsymbol{p}_j, \boldsymbol{p}_j)$$

U, V.

Dot products can be viewed as a **similarity measure**.

$$\boldsymbol{u} \cdot \boldsymbol{v} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \cos \theta$$

$$\boldsymbol{v} = \|\boldsymbol{u}\| \|\boldsymbol{v}\| \cos \theta$$

(c) Gaussian radial basis function (RBF) kernel:

distance  $\Rightarrow$  similarity measure!

$$k_{rbf}(\boldsymbol{p}_i, \boldsymbol{p}_j) = \exp\left(-\frac{\|\boldsymbol{p}_i - \boldsymbol{p}_j\|^2}{2\sigma^2}\right)$$

Calculate  $k_{rbf}(\boldsymbol{p}_i, \boldsymbol{p}_i)$  for all pairs of data points. Explain what the value represents.

р	$k(\boldsymbol{p}_i, \boldsymbol{p}_1)$	$k(\boldsymbol{p}_i, \boldsymbol{p}_2)$	$k(\boldsymbol{p}_i, \boldsymbol{p}_3)$	$k(\boldsymbol{p}_i, \boldsymbol{p}_4)$	$k(\boldsymbol{p}_i, \boldsymbol{p}_5)$
(0,0)	1	$e^{-1}$	$e^{-1}$	$e^{-1}$	$e^{-1}$
(4,4)	$e^{-1}$	1	$e^{-2}$	$e^{-4}$	$e^{-2}$
(-4,4)	$e^{-1}$	$e^{-2}$	1	$e^{-2}$	$e^{-4}$
(-4, -4)	$e^{-1}$	$e^{-4}$	$e^{-2}$	1	$e^{-2}$
(4, -4)	$e^{-1}$	e <sup>-2</sup>	$e^{-4}$	e <sup>-2</sup>	1

р	X	У
1	0	0
2	4	4
3	-4	4
4	-4	-4
5	4	-4

#### Intuition:

- Lines 4-5: Minimize the distortion by reassigning clusters, keeping the centroids unchanged.
- Do not compute the centers (We cannot compute \u03c6(x)...).

Algorithm Kernel K-Means Clustering

- 1: for k = 1 to K do
- 2:  $\mu_k \leftarrow$  random location
- 3: while not converged do
- 4: for i = 1 to m do

5: 
$$c^{(i)} \leftarrow \operatorname{argmin}_k d_{rbf}^2(\boldsymbol{x}^{(i)}, \boldsymbol{\mu}_k)$$

$$\begin{aligned} d_{rbf}^{2}(\boldsymbol{x}^{(i)},\boldsymbol{\mu}_{c}) \\ &= \phi(\boldsymbol{x}^{(i)}) \cdot \phi(\boldsymbol{x}^{(i)}) - 2 \cdot \left(\phi(\boldsymbol{x}^{(i)}) \cdot \frac{1}{\text{size}} \sum_{j \in \mathcal{C}} \phi(\boldsymbol{x}^{(j)})\right) + \left(\left(\frac{1}{\text{size}} \sum_{j \in \mathcal{C}} \phi(\boldsymbol{x}^{(j)})\right) \cdot \left(\frac{1}{\text{size}} \sum_{k \in \mathcal{C}} \phi(\boldsymbol{x}^{(k)})\right)\right) \\ &= k_{rbf}(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(i)}) + \frac{2}{\text{size}} \sum_{j \in \mathcal{C}} k_{rbf}(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)}) + \frac{1}{\text{size}^{2}} \sum_{j,k \in \mathcal{C}} k_{rbf}(\boldsymbol{x}^{(j)}, \boldsymbol{x}^{(k)}) \end{aligned}$$

Extra Slide





Extra Slide

Recap: Principal Component Analysis (PCA)  
$$\boldsymbol{x} = \begin{bmatrix} 4 & 2 & 5 & 1 \\ 1 & 3 & 4 & 0 \end{bmatrix}$$

1. **Mean-center** the data (subtract each  $\mathbf{x}$  by the mean  $\bar{\mathbf{x}}$ ).

$$\bar{\boldsymbol{x}} = \frac{1}{4} \begin{bmatrix} 12\\8 \end{bmatrix} = \begin{bmatrix} 3\\2 \end{bmatrix}$$
$$\hat{\boldsymbol{X}} = \begin{bmatrix} 4 & 2 & 5 & 1\\1 & 3 & 4 & 0 \end{bmatrix} - \begin{bmatrix} 3 & 3 & 3 & 3\\2 & 2 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 1 & -1 & 2 & -2\\-1 & 1 & 2 & -2 \end{bmatrix}$$

2. Compute the covariance matrix  $Cov(\boldsymbol{X}) = \frac{1}{m} \hat{\boldsymbol{X}} \hat{\boldsymbol{X}}^{\top}$ .

$$Cov(\mathbf{X}) = \frac{1}{4} \begin{bmatrix} 1 & -1 & 2 & -2 \\ -1 & 1 & 2 & -2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 2 & 2 \\ -2 & -2 \end{bmatrix} = \begin{bmatrix} 2.5 & 1.5 \\ 1.5 & 2.5 \end{bmatrix}$$



Recap: Principal Component Analysis (PCA)

$$\boldsymbol{X} = \begin{bmatrix} 4 & 2 & 5 & 1 \\ 1 & 3 & 4 & 0 \end{bmatrix}, Cov(\boldsymbol{X}) = \begin{bmatrix} 2.5 & 1.5 \\ 1.5 & 2.5 \end{bmatrix}$$

3. Compute the singuar value decomposition of Cov(X) $(Cov(X) = V\Sigma^2 V^{\top}).$ 

$$\begin{bmatrix} 2.5 & 1.5 \\ 1.5 & 2.5 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

4. Reduce the basis to k components to obtain the new basis  $\tilde{\boldsymbol{U}}$ . Calculate the ratio of explained variance  $\frac{\sum_{i=1}^{r} \sigma_i^2}{\sum_{i=1}^{m} \sigma_i^2}$ .

$$\tilde{\boldsymbol{U}} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}, \frac{\sum_{i=1}^{r} \sigma_i^2}{\sum_{i=1}^{m} \sigma_i^2} = \frac{4}{4+1} = 80\%$$



Recap: Principal Component Analysis (PCA)

$$oldsymbol{X} = egin{bmatrix} 4 & 2 & 5 & 1 \ 1 & 3 & 4 & 0 \end{bmatrix}, oldsymbol{ ilde{U}} = egin{bmatrix} rac{1}{\sqrt{2}} \ rac{1}{\sqrt{2}} \end{bmatrix}$$

**•** Reduction:  $\boldsymbol{Z} = \boldsymbol{\tilde{\boldsymbol{U}}}^{\top} \boldsymbol{X}$ 

$$Z = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 4 & 2 & 5 & 1 \\ 1 & 3 & 4 & 0 \end{bmatrix} = \begin{bmatrix} \frac{5}{\sqrt{2}} & \frac{5}{\sqrt{2}} & \frac{9}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

• Reconstruction:  $\boldsymbol{X} \approx \tilde{\boldsymbol{U}} \boldsymbol{Z}$ 

$$X \approx \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{5}{\sqrt{2}} & \frac{5}{\sqrt{2}} & \frac{9}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 2.5 & 2.5 & 4.5 & 0.5 \\ 2.5 & 2.5 & 4.5 & 0.5 \end{bmatrix}$$



We are only required to store  $\tilde{U}$  and Z (instead of X)!

(a) The current choice of k = 9 does not produce a very nice output. What is a good value for k? Justify your answer.





After k=9

- 1. Compute the covariance matrix  $Cov(\mathbf{X}) = \frac{1}{m} \hat{\mathbf{X}} \hat{\mathbf{X}}^{\top}$  (already includes mean-centering).
- 2. Compute the singuar value decomposition of Cov(X) ( $Cov(X) = V\Sigma^2 V^{\top}$ ).

```
1 cov = np.cov(X, bias=True)
2 U, s, VT = np.linalg.svd(cov)
```

Goal: Retain at least 99% of the explained variance.

```
def explained_variance(s, k):
      """ returns the explained variance when we keep the first k cols
                                                                              0.0.0
2
      return sum(s[:k]) / sum(s)
3
4
5 minimum_retained_variance = 0.99
6
  k = 0
7 \text{ var} = 0
8 while var < minimum_retained_variance:</pre>
      k += 1
9
  var = explained_variance(s, k)
10
```

Answer: k = 286



After



- (b) For the value of k you select in (a), what is the space saved by doing this compression?
- Original: X (512 × 1536)
  New: Ũ (512 × 286) and Z (286 × 1536).
  Compression ratio = (512 × 286) + (286 × 1536) 512 × 1536 = 0.745.
  25.5% space saved!

Reduction:  $\boldsymbol{Z} = \boldsymbol{\tilde{U}}^{\top} \boldsymbol{X}$ Reconstruction:  $\boldsymbol{X} \approx \boldsymbol{\tilde{U}} \boldsymbol{Z}$ 

- (c) What are the drawbacks of this form of compression?
- ▶ This is a lossy compression as we do not regain 100% of the variance.
- *Ũ* is only used for one image ⇒ waste of space (if k is large, we might end up using more space).

Refer to the template code provided.

- 1. Let  $\mathcal{K}^{(1)}(\boldsymbol{u}, \boldsymbol{v})$  and  $\mathcal{K}^{(2)}(\boldsymbol{u}, \boldsymbol{v})$  be kernel functions. Work out the transformed features in the kernel  $\mathcal{K}^{(1)}(\boldsymbol{u}, \boldsymbol{v}) + \mathcal{K}^{(2)}(\boldsymbol{u}, \boldsymbol{v})$ .
- Let K<sup>(1)</sup>(u, v) and K<sup>(2)</sup>(u, v) be kernel functions.
   Work out the transformed features in the kernel K<sup>(1)</sup>(u, v) × K<sup>(2)</sup>(u, v).
- 3. Invent your own kernel!

Extra Slide

## Bonus. Composing Kernels

#### Sample Solution:

1. We have

$$egin{aligned} &\mathcal{K}^{(1)}(oldsymbol{u},oldsymbol{v}) = \phi^{(1)}(oldsymbol{u}) \cdot \phi_1(oldsymbol{v}) \ &\mathcal{K}^{(2)}(oldsymbol{u},oldsymbol{v}) = \phi^{(2)}(oldsymbol{u}) \cdot \phi_2(oldsymbol{v}) \end{aligned}$$

Then, we get

$$\mathcal{K}^{(1)}(\boldsymbol{u},\boldsymbol{v}) + \mathcal{K}^{(2)}(\boldsymbol{u},\boldsymbol{v}) = \begin{bmatrix} \phi^{(1)}(\boldsymbol{u}) \\ \phi^{(2)}(\boldsymbol{u}) \end{bmatrix} \cdot \begin{bmatrix} \phi^{(1)}(\boldsymbol{v}) \\ \phi^{(2)}(\boldsymbol{v}) \end{bmatrix}$$

which means the addition of two kernel functions essentially **concatenates** the transformed features together.

Extra Slide



### Bonus. Composing Kernels

2. Similarly, we get

þ

$$\begin{split} \mathcal{K}^{(1)}(\boldsymbol{u},\boldsymbol{v}) \times \mathcal{K}^{(2)}(\boldsymbol{u},\boldsymbol{v}) &= (\phi^{(1)}(\boldsymbol{u}) \cdot \phi^{(1)}(\boldsymbol{v})) \times (\phi^{(2)}(\boldsymbol{u}) \cdot \phi^{(2)}(\boldsymbol{v})) \\ &= \left(\sum_{i} \phi^{(1)}_{i}(\boldsymbol{u}) \phi^{(1)}_{i}(\boldsymbol{v})\right) \times \left(\sum_{j} \phi^{(2)}_{j}(\boldsymbol{u}) \phi^{(2)}_{j}(\boldsymbol{v})\right) \\ &= \sum_{ij} \left(\phi^{(1)}_{i}(\boldsymbol{u}) \phi^{(2)}_{j}(\boldsymbol{u})\right) \left(\phi^{(1)}_{i}(\boldsymbol{v}) \phi^{(2)}_{j}(\boldsymbol{v})\right) \end{split}$$

Define another set of transformed features  $\phi^*(\boldsymbol{u})$  which consists of  $\phi_i^{(1)}(\boldsymbol{u})\phi_j^{(2)}(\boldsymbol{u})$  for all pairs (i,j), then

$$\mathcal{K}^{(1)}(\boldsymbol{u},\boldsymbol{v}) imes \mathcal{K}^{(2)}(\boldsymbol{u},\boldsymbol{v}) = \phi^*(\boldsymbol{u}) \cdot \phi^*(\boldsymbol{v})$$

which means the product of two kernel functions takes the **pairwise product** of between their transformed features.

## Bonus. Composing Kernels

#### 1. Just a simple concatenation.

```
1 def transform_add(X1, X2):
2 return np.hstack([X1, X2])
```

2. Broadcasting is your best friend!

```
1 def transform_multiply(X1, X2):
2 n = X1.shape[0]
3 pairwise_prod = X1.reshape(n, -1, 1) * X2.reshape(n, 1, -1)
4 return pairwise_prod.reshape(n, -1)
```