# Introduction to Kernel Methods for Comparing Distributions

Wittawat Jitkrittum

Max Planck Institute for Intelligent Systems, Germany wittawatj@gmail.com

> Bangkok Machine Learning Meetup 7 March 2018

Have: Two collections of samples X, Y from unknown distributions p and q.

Two goals: using only  $X, Y \dots$ 

1 Measure the distance between p and q.

2 Are p and q different (not just by chance)?

 $\implies$  Two-sample testing

Have: Two collections of samples X, Y from unknown distributions p and q.

### Two goals: using only $X, Y \dots$

- 1 Measure the distance between p and q.
- 2 Are p and q different (not just by chance)?
  - $\implies$  Two-sample testing

Have: Two collections of samples X, Y from unknown distributions p and q.

Two goals: using only  $X, Y \dots$ 

1 Measure the distance between p and q.

2 Are p and q different (not just by chance)?  $\implies$  Two-sample testing



Have: Two collections of samples X, Y from unknown distributions p and q.

Two goals: using only  $X, Y \dots$ 

- 1 Measure the distance between p and q.
- 2 Are p and q different (not just by chance)?
  - $\implies$  Two-sample testing







Create a Meetup

Explore Messages Notifications

**()** ~

Wednesday, March 7, 2018

#### "Meta-Learning" and "Kernel Methods for Comparing Distribution"



Hosted by James, Sorawit and 3 others From BKK Machine Learning

# Ads at location 1

#### Details

Bangkok Machine Learning meetup is back!

This time we will learn about two interesting topics in machine learning from two very special guests, Sam Witterveen and Wittawat Jitkrittum.

Sam is a Google Developer Expert in Machine Learning. He regularly shares his knowledge at events and trainings across Asia and is coorganiser of the Singapore TensorFlow and Deep Learning group Wednesday, March 7, 2018 6:00 PM to 8:30 PM Add to calendar

Room 206, 2nd Floor, 50 Years Anusorn Building (BBA Building), Chuilalongkorn Business School อาคารอนุสรณ์ 50ปี -Bangkok Close to MRT Sam Yan Station



#### Details

Bangkok Machine Learning meetup is back!

This time we will learn about two interesting topics in machine learning from two very special guests, Sam Witterveen and Wittawat Jitkrittum.

Sam is a Google Developer Expert in Machine Learning. He regularly shares his knowledge at events and trainings across Asia and is coorganiser of the Singapore TensorFlow and Deep Learning group

# Ads at location 2

อาคารอนุสรณ์ 50ปี Bangkok Close to MRT Sam Yan Station

Does ads at location 1 have the same effect as ads at location 2?

Does ads at location 1 have the same effect as ads at location 2?

- X = time it takes users to click ads at location 1.
- Y =time it takes users to click ads at location 2.

Does ads at location 1 have the same effect as ads at location 2?

- X = time it takes users to click ads at location 1.
- Y =time it takes users to click ads at location 2.



# Application 2: Data Integration



Data collected from lab 1: X.

Data collected from lab 2: Y.

■ If they have different distributions, do not merge.

# Application 3: Benchmarking Generative Models





Observed MNIST handwritten digits. X.

Generated images from a model. Y.

Is Y similar to X?

Distance between distributions can be used to train generative models.



#### 1 Background

2 Kernel Methods for Comparing Distributions

3 Nonparametric Two-Sample Testing

4 Further Topics and Conclusion



#### 1 Background

2 Kernel Methods for Comparing Distributions

3 Nonparametric Two-Sample Testing

4 Further Topics and Conclusion

Let a = (a<sub>1</sub>,..., a<sub>d</sub>)<sup>T</sup>, b and c be vectors in ℝ<sup>d</sup>.
Norm (length): ||a|| := √∑<sub>i=1</sub><sup>d</sup> a<sub>i</sub><sup>2</sup>.

Inner product (dot product)

$$egin{array}{ll} m{a}\cdotm{b} &= m{a}^{ op}m{b} &= \langlem{a},m{b}
angle = \sum_{i=1}^d a_ib_i \ &= a_1b_1 + a_2b_2 + \cdots + a_db_d. \end{array}$$

 $\bullet \langle a, b \rangle = \text{similarity between } a \text{ and } b.$ 

An inner product induces a norm:  $||a|| = \sqrt{\langle a, a \rangle}$ .

Inner product (dot product)

$$egin{array}{ll} m{a}\cdotm{b} &= m{a}^{ op}m{b} &= \langlem{a},m{b}
angle &= \sum\limits_{i=1}^d a_ib_i \ &= a_1b_1 + a_2b_2 + \cdots + a_db_d. \end{array}$$

 $\blacksquare \langle a, b \rangle = \text{similarity between } a \text{ and } b.$ 

• An inner product induces a norm:  $||a|| = \sqrt{\langle a, a \rangle}$ .

Inner product (dot product)

1

b

⟨a, b⟩ = similarity between a and b.
An inner product induces a norm: ||a|| = √⟨a, a⟩.

Inner product (dot product)

$$a \cdot b = a^{\top}b = \langle a, b \rangle = \sum_{i=1}^{d} a_i b_i$$
  
=  $a_1 b_1 + a_2 b_2 + \dots + a_d b_d$ .  
 $a \cdot b$ 

 $\bullet \langle a, b \rangle = \text{similarity between } a \text{ and } b.$ 

• An inner product induces a norm:  $||a|| = \sqrt{\langle a, a \rangle}$ .

Three properties:

- 1 (Linear):  $\langle \alpha a + \beta b, c \rangle = \alpha \langle a, c \rangle + \beta \langle b, c \rangle$
- 2 (Symmetric):  $\langle \boldsymbol{a}, \boldsymbol{b} \rangle = \langle \boldsymbol{b}, \boldsymbol{a} \rangle$
- 3  $\langle a, a \rangle \geq 0$  and  $\langle a, a \rangle = 0$  if and only if a = 0.
- For x, y ∈ ℝ, we have (x − y)<sup>2</sup> = x<sup>2</sup> − 2xy + y<sup>2</sup>.
  In general: ||a − b||<sup>2</sup> = ⟨a, a⟩ − 2⟨a, b⟩ + ⟨b, b⟩.
  ||a − b|| = distance between a and b.

#### Definition 1 (Hilbert space).

A Hilbert space  $\mathcal{H}$  is a complete inner product space.

- Hilbert space  $\approx$  a space with an inner product defined.  $\mathbb{R}^d$  is a Hilbert space.
- In general, can be a space of generic objects.

Three properties:

- 1 (Linear):  $\langle \alpha a + \beta b, c \rangle = \alpha \langle a, c \rangle + \beta \langle b, c \rangle$
- 2 (Symmetric):  $\langle \boldsymbol{a}, \boldsymbol{b} \rangle = \langle \boldsymbol{b}, \boldsymbol{a} \rangle$
- 3  $\langle a, a \rangle \geq 0$  and  $\langle a, a \rangle = 0$  if and only if a = 0.
- For x, y ∈ ℝ, we have (x − y)<sup>2</sup> = x<sup>2</sup> − 2xy + y<sup>2</sup>.
  In general: ||a − b||<sup>2</sup> = ⟨a, a⟩ − 2 ⟨a, b⟩ + ⟨b, b⟩.
  ||a − b|| = distance between a and b.

#### Definition 1 (Hilbert space).

A Hilbert space  $\mathcal{H}$  is a complete inner product space.

- Hilbert space  $\approx$  a space with an inner product defined.  $\mathbb{R}^d$  is a Hilbert space.
- In general, can be a space of generic objects.

Three properties:

- 1 (Linear):  $\langle \alpha a + \beta b, c \rangle = \alpha \langle a, c \rangle + \beta \langle b, c \rangle$
- 2 (Symmetric):  $\langle \boldsymbol{a}, \boldsymbol{b} \rangle = \langle \boldsymbol{b}, \boldsymbol{a} \rangle$
- 3  $\langle a, a \rangle \geq 0$  and  $\langle a, a \rangle = 0$  if and only if a = 0.
- For x, y ∈ ℝ, we have (x − y)<sup>2</sup> = x<sup>2</sup> − 2xy + y<sup>2</sup>.
  In general: ||a − b||<sup>2</sup> = ⟨a, a⟩ − 2 ⟨a, b⟩ + ⟨b, b⟩.
  ||a − b|| = distance between a and b.

#### Definition 1 (Hilbert space).

- A Hilbert space  $\mathcal{H}$  is a complete inner product space.
  - Hilbert space  $\approx$  a space with an inner product defined.  $\mathbb{R}^d$  is a Hilbert space.
  - In general, can be a space of generic objects.

Three properties:

- 1 (Linear):  $\langle \alpha a + \beta b, c \rangle = \alpha \langle a, c \rangle + \beta \langle b, c \rangle$
- 2 (Symmetric):  $\langle \boldsymbol{a}, \boldsymbol{b} \rangle = \langle \boldsymbol{b}, \boldsymbol{a} \rangle$
- 3  $\langle a, a \rangle \geq 0$  and  $\langle a, a \rangle = 0$  if and only if a = 0.
- For x, y ∈ ℝ, we have (x − y)<sup>2</sup> = x<sup>2</sup> − 2xy + y<sup>2</sup>.
  In general: ||a − b||<sup>2</sup> = ⟨a, a⟩ − 2 ⟨a, b⟩ + ⟨b, b⟩.
  ||a − b|| = distance between a and b.

#### Definition 1 (Hilbert space).

A Hilbert space  $\mathcal{H}$  is a complete inner product space.

- Hilbert space  $\approx$  a space with an inner product defined.  $\mathbb{R}^d$  is a Hilbert space.
- In general, can be a space of generic objects.



#### 1 Background

#### 2 Kernel Methods for Comparing Distributions

#### 3 Nonparametric Two-Sample Testing

#### 4 Further Topics and Conclusion



Two Gaussian distributions.





We have only samples X ~ p and Y ~ q.
X = {x<sub>1</sub>,..., x<sub>n</sub>} and Y = {y<sub>1</sub>,..., y<sub>n</sub>}. Sets of numbers.



We have only samples X ~ p and Y ~ q.
X = {x<sub>1</sub>,..., x<sub>n</sub>} and Y = {y<sub>1</sub>,..., y<sub>n</sub>}. Sets of numbers.



Assume no differece in high-order moments.
"Distance" = difference in the means. T-test.

 $\begin{array}{l} \text{(population)} \ D_1(p,q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]| \\ \text{(empirical)} \ \hat{D}_1(\mathsf{X},\mathsf{Y}) = \left| \frac{1}{n}\sum_{i=1}^n x_i - \frac{1}{n}\sum_{j=1}^n y_j \right| \end{array}$ 



•  $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?

**Idea:** look at difference in means of features  $\phi(\cdot)$  of X and Y.



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?
- **Idea**: look at difference in means of features  $\phi(\cdot)$  of X and Y.



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?
- **Idea**: look at difference in means of features  $\phi(\cdot)$  of X and Y.



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?
- **Idea**: look at difference in means of features  $\phi(\cdot)$  of X and Y.



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?
- **Idea**: look at difference in means of features  $\phi(\cdot)$  of X and Y.
- New "distance":

$$D_2(p, q) = ig\|\mathbb{E}_{X \sim p}[\phi(X)] - \mathbb{E}_{Y \sim q}[\phi(X)]ig\|,$$
  
where  $\phi(x) = (x, x^2)^ op$ .



•  $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference. Why?

**Idea**: look at difference in means of features  $\phi(\cdot)$  of X and Y.

$$D_2(p, q) = \left\| \left( egin{array}{c} \mathbb{E}_{X \sim p}[X] \ \mathbb{E}_{X \sim p}[X^2] \end{array} 
ight) - \left( egin{array}{c} \mathbb{E}_{X \sim q}[X] \ \mathbb{E}_{X \sim q}[X^2] \end{array} 
ight) 
ight\|$$

# Case 3: Difference in High-Order Moments



\$\phi(x) = (x, x^2, x^4)^\texts\$ works. Difference is in kurtosis (4<sup>th</sup> moment).
\$\phi(x) = (x, x^2, x^4, \cos x, e^x, ...)^\texts\$. But, when to stop?

Solution: Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

# Case 3: Difference in High-Order Moments



φ(x) = (x, x<sup>2</sup>, x<sup>4</sup>)<sup>T</sup> works. Difference is in kurtosis (4<sup>th</sup> moment).
φ(x) = (x, x<sup>2</sup>, x<sup>4</sup>, cos x, e<sup>x</sup>, ...)<sup>T</sup>. But, when to stop?
Solution: Use an infinite-dimensional feature map φ(·) with the kernel trick.
## Case 3: Difference in High-Order Moments



- \$\phi(x) = (x, x^2, x^4)^\top works. Difference is in kurtosis (4<sup>th</sup> moment).
  \$\phi(x) = (x, x^2, x^4, \cos x, e^x, ...)^\top.\$ But, when to stop?
- **Solution:** Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

## Case 3: Difference in High-Order Moments



- \$\phi(x) = (x, x^2, x^4)^\top works. Difference is in kurtosis (4<sup>th</sup> moment).
  \$\phi(x) = (x, x^2, x^4, \cos x, e^x, ...)^\top.\$ But, when to stop?
- **Solution:** Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

# The (Kernel) Mean Embedding [Smola et al., 2007]

- Given a feature map  $\phi(\cdot)$  mapping to a Hilbert space  $\mathcal{H}$ ,
  - represent p with  $\mu_p := \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})]$  (i.e., the mean embedding of p),
  - represent q with  $\mu_q := \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})].$

## The (Kernel) Mean Embedding [Smola et al., 2007]

- Given a feature map  $\phi(\cdot)$  mapping to a Hilbert space  $\mathcal{H}$ ,
  - represent p with  $\mu_p := \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})]$  (i.e., the mean embedding of p),
  - represent q with  $\mu_q := \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})].$

## The (Kernel) Mean Embedding [Smola et al., 2007]

- Given a feature map  $\phi(\cdot)$  mapping to a Hilbert space  $\mathcal{H}$ ,
  - represent p with  $\mu_p := \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})]$  (i.e., the mean embedding of p),
  - represent q with  $\mu_q := \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})].$



**\square**  $\mathcal{H}$  can be infinite dimensional. Depends on  $\phi(\cdot)$ .

- If  $\phi(x) = (x, x^2)^{ op}$ , then  $\mathcal{H} = \mathbb{R}^2$ .
- Then, measure the distance in  $\mathcal{H}$ .
- The distance is called the "Maximum Mean Discrepancy" (MMD).

14/44

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}}$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ . Recall  $\|\boldsymbol{a} - \boldsymbol{b}\|^2 = \langle \boldsymbol{a}, \boldsymbol{a} \rangle - 2 \langle \boldsymbol{a}, \boldsymbol{b} \rangle + \langle \boldsymbol{b}, \boldsymbol{b} \rangle.$ 

Depend on only the inner product (φ(x), φ(y)).
 Don't need φ(x) explicitly (could be ∞-dimensional!).

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ . • Recall  $||a - b||^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .  $MMD^2(p, q) = \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2$   $= (\mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x \sim p}[\phi(x')]) - 2 (\mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)])$   $+ (\mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y \sim q}[\phi(y')])$   $= \mathbb{E}_{x \sim p}\mathbb{E}_{x' \sim p}(\phi(x), \phi(x')) - 2\mathbb{E}_{x \sim p}\mathbb{E}_{y \sim q}(\phi(x), \phi(y))$  $+ \mathbb{E}_{y \sim q}\mathbb{E}_{y' \sim q}(\phi(y), \phi(y'))$ 

Depend on only the inner product ⟨φ(x), φ(y)⟩.
 Don't need φ(x) explicitly (could be ∞-dimensional!).

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

• Recall  $||a - b||^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$ext{MMD}^2(p, oldsymbol{q}) = \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim oldsymbol{q}}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}}^2$$

$$\begin{split} &= \langle \mathbb{E}_{\boldsymbol{x} \sim \boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{x}' \sim \boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x}')] \rangle - 2 \left\langle \mathbb{E}_{\boldsymbol{x} \sim \boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})] \right\rangle \\ &+ \left\langle \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})], \mathbb{E}_{\boldsymbol{y}' \sim \boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y}')] \right\rangle \end{split}$$

 $egin{aligned} &= \mathbb{E}_{oldsymbol{x} \sim oldsymbol{p}} \mathbb{E}_{oldsymbol{x}' \sim oldsymbol{p}} ig \langle \phi(oldsymbol{x}), \phi(oldsymbol{x}') 
angle - 2\mathbb{E}_{oldsymbol{x} \sim oldsymbol{p}} \mathbb{E}_{oldsymbol{y} \sim oldsymbol{q}} ig \langle \phi(oldsymbol{x}), \phi(oldsymbol{y}') 
angle \ &+ \mathbb{E}_{oldsymbol{y} \sim oldsymbol{q}} \mathbb{E}_{oldsymbol{y}' \sim oldsymbol{q}} ig \langle \phi(oldsymbol{y}), \phi(oldsymbol{y}') 
angle \end{aligned}$ 

Depend on only the inner product ⟨φ(x), φ(y)⟩.
Don't need φ(x) explicitly (could be ∞-dimensional!).

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

**Recall**  $\|\boldsymbol{a} - \boldsymbol{b}\|^2 = \langle \boldsymbol{a}, \boldsymbol{a} \rangle - 2 \langle \boldsymbol{a}, \boldsymbol{b} \rangle + \langle \boldsymbol{b}, \boldsymbol{b} \rangle.$ 

$$\begin{split} \mathsf{MMD}^2(\boldsymbol{p},\boldsymbol{q}) &= \left\| \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{x}'\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x}')] \rangle - 2 \left\langle \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})] \right\rangle \\ &+ \left\langle \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})], \mathbb{E}_{\boldsymbol{y}'\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y}')] \right\rangle \\ &= \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}} \mathbb{E}_{\boldsymbol{x}'\sim\boldsymbol{p}} \langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}') \rangle - 2\mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}} \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}} \langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{y}) \rangle \end{split}$$

 $+ \mathbb{E}_{oldsymbol{y} \sim oldsymbol{q}} \mathbb{E}_{oldsymbol{y}' \sim oldsymbol{q}} ig\langle \phi(oldsymbol{y}), \phi(oldsymbol{y}') ig
angle$ 

Depend on only the inner product ⟨φ(x), φ(y)⟩.
Don't need φ(x) explicitly (could be ∞-dimensional!).

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

**Recall**  $\|\boldsymbol{a} - \boldsymbol{b}\|^2 = \langle \boldsymbol{a}, \boldsymbol{a} \rangle - 2 \langle \boldsymbol{a}, \boldsymbol{b} \rangle + \langle \boldsymbol{b}, \boldsymbol{b} \rangle.$ 

$$\begin{split} \mathsf{MMD}^2(\boldsymbol{p},\boldsymbol{q}) &= \left\| \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{x}'\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x}')] \rangle - 2 \langle \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}}[\boldsymbol{\phi}(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})] \rangle \\ &+ \langle \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y})], \mathbb{E}_{\boldsymbol{y}'\sim\boldsymbol{q}}[\boldsymbol{\phi}(\boldsymbol{y}')] \rangle \\ &= \mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}} \mathbb{E}_{\boldsymbol{x}'\sim\boldsymbol{p}} \langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}') \rangle - 2\mathbb{E}_{\boldsymbol{x}\sim\boldsymbol{p}} \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}} \langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{y}) \rangle \\ &+ \mathbb{E}_{\boldsymbol{y}\sim\boldsymbol{q}} \mathbb{E}_{\boldsymbol{y}'\sim\boldsymbol{q}} \langle \boldsymbol{\phi}(\boldsymbol{y}), \boldsymbol{\phi}(\boldsymbol{y}') \rangle \end{split}$$

Depend on only the inner product ⟨φ(x), φ(y)⟩.
Don't need φ(x) explicitly (could be ∞-dimensional!).

$$ext{MMD}(p,q) := \left\| \mathbb{E}_{oldsymbol{x} \sim p}[\phi(oldsymbol{x})] - \mathbb{E}_{oldsymbol{y} \sim q}[\phi(oldsymbol{y})] 
ight\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

**Recall**  $\|\boldsymbol{a} - \boldsymbol{b}\|^2 = \langle \boldsymbol{a}, \boldsymbol{a} \rangle - 2 \langle \boldsymbol{a}, \boldsymbol{b} \rangle + \langle \boldsymbol{b}, \boldsymbol{b} \rangle.$ 

$$\begin{split} \mathrm{MMD}^2(p, \boldsymbol{q}) &= \left\| \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{x}' \sim p}[\phi(\boldsymbol{x}')] \rangle - 2 \left\langle \mathbb{E}_{\boldsymbol{x} \sim p}[\phi(\boldsymbol{x})], \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})] \right\rangle \\ &+ \left\langle \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}}[\phi(\boldsymbol{y})], \mathbb{E}_{\boldsymbol{y}' \sim \boldsymbol{q}}[\phi(\boldsymbol{y}')] \right\rangle \\ &= \mathbb{E}_{\boldsymbol{x} \sim p} \mathbb{E}_{\boldsymbol{x}' \sim p} \left\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \right\rangle - 2\mathbb{E}_{\boldsymbol{x} \sim p} \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}} \left\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{y}) \right\rangle \\ &+ \mathbb{E}_{\boldsymbol{y} \sim \boldsymbol{q}} \mathbb{E}_{\boldsymbol{y}' \sim \boldsymbol{q}} \left\langle \phi(\boldsymbol{y}), \phi(\boldsymbol{y}') \right\rangle \end{split}$$

Depend on only the inner product (φ(x), φ(y)).
Don't need φ(x) explicitly (could be ∞-dimensional!).

 $\begin{array}{l} \blacksquare \text{ Define } k(x,x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}} \text{ (kernel).} \\ \\ \text{MMD}_{k}^{2}(p,q) = \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^{2} \\ \\ = \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} k(x,x') - 2\mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} k(x,y') \\ \\ + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} k(y,y'). \end{array}$ 

Unbiased estimator:

$$egin{aligned} \widehat{ ext{MMD}}_k^2(\mathsf{X},\mathsf{Y}) &= rac{1}{n(n-1)}\sum_{i=1}^n\sum_{j
eq i}k(x_i,x_j) - rac{2}{n^2}\sum_{i=1}^n\sum_{j=1}^nk(x_i,y_j) \ &+ rac{1}{n(n-1)}\sum_{i=1}^n\sum_{j
eq i}k(y_i,y_j). \end{aligned}$$

 $\mathbf{z}$   $k(x,x') pprox ext{similarity between } x ext{ and } x'.$ 

$$\begin{array}{l} \bullet \quad \text{Define } k(x,x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}} \text{ (kernel).} \\ \\ & \text{MMD}_{k}^{2}(p,q) = \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^{2} \\ \\ & = \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} k(x,x') - 2\mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} k(x,y') \\ \\ & + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} k(y,y'). \end{array}$$

Unbiased estimator:

$$egin{aligned} \widehat{ ext{MMD}}_k^2({\sf X},{\sf Y}) &= rac{1}{n(n-1)}\sum_{i=1}^n\sum_{j
eq i}k(x_i,x_j) - rac{2}{n^2}\sum_{i=1}^n\sum_{j=1}^nk(x_i,y_j) \ &+ rac{1}{n(n-1)}\sum_{i=1}^n\sum_{j
eq i}k(y_i,y_j). \end{aligned}$$

•  $k(x, x') \approx ext{similarity between } x ext{ and } x'.$ 

# Intuition for the MMD

**Dogs**  $\sim p$  and fish  $\sim q$ .

Each entry is one of  $k(\text{dog}_i, \text{dog}_j)$ ,  $k(\text{dog}_i, \text{fish}_j)$ , or  $k(\text{fish}_i, \text{fish}_j)$ 



Intuition for the MMD



#### • Defining k(x, x') from $\phi(\cdot)$ is always valid.

Can start directly from k(x, x') without specifying φ(·).
What k is valid?

Definition 2 (Positive definite kernel).

A <u>symmetric</u> function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is called <u>positive definite</u> if, for any integer  $n > 0, c_1, \ldots, c_n \in \mathbb{R}$ , and  $x_1, \ldots, x_n \in \mathcal{X}$ , we have  $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \ge 0.$ 

Equivalently, the Gram matrix K is a positive semi-definite matrix where  $(K)_{ij} = k(x_i, x_j)$ .

• x can be anything (e.g., vector, image, tree, string, graph, ...).

- Defining k(x, x') from  $\phi(\cdot)$  is always valid.
- Can start directly from k(x, x') without specifying  $\phi(\cdot)$ .
- What *k* is valid?

Definition 2 (Positive definite kernel).

A <u>symmetric</u> function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is called <u>positive definite</u> if, for any integer  $n > 0, c_1, \ldots, c_n \in \mathbb{R}$ , and  $x_1, \ldots, x_n \in \mathcal{X}$ , we have  $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \ge 0.$ 

- Equivalently, the Gram matrix K is a positive semi-definite matrix where  $(K)_{ij} = k(x_i, x_j)$ .
- x can be anything (e.g., vector, image, tree, string, graph, ...).

- Defining k(x, x') from  $\phi(\cdot)$  is always valid.
- Can start directly from k(x, x') without specifying  $\phi(\cdot)$ .
- What k is valid?

### Definition 2 (Positive definite kernel).

A <u>symmetric</u> function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is called <u>positive definite</u> if, for any integer  $n > 0, c_1, \ldots, c_n \in \mathbb{R}$ , and  $x_1, \ldots, x_n \in \mathcal{X}$ , we have  $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \ge 0.$ 

• Equivalently, the Gram matrix K is a positive semi-definite matrix where  $(K)_{ij} = k(x_i, x_j)$ .

**x** can be *anything* (e.g., vector, image, tree, string, graph, ...).

- Defining k(x, x') from  $\phi(\cdot)$  is always valid.
- Can start directly from k(x, x') without specifying  $\phi(\cdot)$ .
- What k is valid?

### Definition 2 (Positive definite kernel).

A <u>symmetric</u> function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is called <u>positive definite</u> if, for any integer  $n > 0, c_1, \ldots, c_n \in \mathbb{R}$ , and  $x_1, \ldots, x_n \in \mathcal{X}$ , we have  $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \ge 0.$ 

- Equivalently, the Gram matrix K is a positive semi-definite matrix where  $(K)_{ij} = k(x_i, x_j)$ .
- x can be anything (e.g., vector, image, tree, string, graph, ...).

#### Theorem 1 (Moore-Aronszajn).

Assume  $k(\cdot, \cdot)$  is positive definite.

- **1** k is an inner product in some Hilbert space  $\mathcal{H}$ .
- 2 There exists  $\phi(\cdot)$  such that  $k(x,y) = \langle \phi(x), \phi(y) 
  angle_{\mathcal{H}}.$

Summary: Pos. def. k automatically defines  $\phi(\cdot)$  (implicitly).

- Defining k can be easier than defining  $\phi(\cdot)$ . Imagine strings.
- To study  $\mathcal{H}$ , can study  $k(\cdot, \cdot)$  instead of  $\phi(\cdot)$ .
  - Reproducing kernel Hilbert spaces (RKHS).

#### Theorem 1 (Moore-Aronszajn).

Assume  $k(\cdot, \cdot)$  is positive definite.

- **1** k is an inner product in some Hilbert space  $\mathcal{H}$ .
- 2 There exists  $\phi(\cdot)$  such that  $k(x, y) = \langle \phi(x), \phi(y) 
  angle_{\mathcal{H}}$ .

Summary: Pos. def. k automatically defines  $\phi(\cdot)$  (implicitly).

- Defining k can be easier than defining  $\phi(\cdot)$ . Imagine strings.
- To study  $\mathcal{H}$ , can study  $k(\cdot, \cdot)$  instead of  $\phi(\cdot)$ .
  - Reproducing kernel Hilbert spaces (RKHS).

#### Theorem 1 (Moore-Aronszajn).

Assume  $k(\cdot, \cdot)$  is positive definite.

- **1** k is an inner product in some Hilbert space  $\mathcal{H}$ .
- 2 There exists  $\phi(\cdot)$  such that  $k(x, y) = \langle \phi(x), \phi(y) 
  angle_{\mathcal{H}}$ .

Summary: Pos. def. k automatically defines  $\phi(\cdot)$  (implicitly).

Defining k can be easier than defining  $\phi(\cdot)$ . Imagine strings.

To study  $\mathcal{H}$ , can study  $k(\cdot, \cdot)$  instead of  $\phi(\cdot)$ .

• Reproducing kernel Hilbert spaces (RKHS).

#### Theorem 1 (Moore-Aronszajn).

Assume  $k(\cdot, \cdot)$  is positive definite.

- **1** k is an inner product in some Hilbert space  $\mathcal{H}$ .
- 2 There exists  $\phi(\cdot)$  such that  $k(x, y) = \langle \phi(x), \phi(y) 
  angle_{\mathcal{H}}$ .

Summary: Pos. def. k automatically defines  $\phi(\cdot)$  (implicitly).

- Defining k can be easier than defining  $\phi(\cdot)$ . Imagine strings.
- To study  $\mathcal{H}$ , can study  $k(\cdot, \cdot)$  instead of  $\phi(\cdot)$ .
  - Reproducing kernel Hilbert spaces (RKHS).

Let  $\mathcal{X} = \mathbb{R}^d$  (domain).

 $k({m x},{m y}) = \left({m x}^{ op}{m y}
ight)^m$ 

is positive definite, for  $m \in \{1, 2, \ldots\}$ .

Consider 
$$d = 2$$
 and  $m = 3$ .  
 $k\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}\right) = (x_1y_1 + x_2y_2)^3.$ 

Feature map

$$\phi\left(\left(egin{array}{c} x_1\ x_2\end{array}
ight)
ight)=\left(egin{array}{c} x_1^3\ \sqrt{3}x_1^2x_2\ \sqrt{3}x_1x_2^2\ x_2^3\end{array}
ight)=\left(egin{array}{c} \phi_1(x)\ \phi_2(x)\ \phi_2(x)\ \phi_3(x)\ \phi_4(x)\end{array}
ight)$$

So,  $\mathcal{H} = \mathbb{R}^4$ . Show that  $k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \phi(x)^\top \phi(y)$ .

Let  $\mathcal{X} = \mathbb{R}^d$  (domain).

$$k(oldsymbol{x},oldsymbol{y}) = \left(oldsymbol{x}^{ op}oldsymbol{y}
ight)^m$$

is positive definite, for  $m \in \{1, 2, \ldots\}$ .

• Consider 
$$d = 2$$
 and  $m = 3$ .  
 $k\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}\right) = (x_1y_1 + x_2y_2)^3.$ 

Feature map

$$\phi\left(\left(egin{array}{c} x_1\ x_2\end{array}
ight)
ight)=\left(egin{array}{c} x_1^3\ \sqrt{3}x_1^2x_2\ \sqrt{3}x_1x_2^2\ x_2^3\end{array}
ight)=\left(egin{array}{c} \phi_1(x)\ \phi_2(x)\ \phi_2(x)\ \phi_3(x)\ \phi_4(x)\end{array}
ight)$$

So,  $\mathcal{H} = \mathbb{R}^4$ . Show that  $k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \phi(x)^\top \phi(y)$ .

Let  $\mathcal{X} = \mathbb{R}^d$  (domain).

$$k(oldsymbol{x},oldsymbol{y}) = \left(oldsymbol{x}^{ op}oldsymbol{y}
ight)^m$$

is positive definite, for  $m \in \{1, 2, \ldots\}$ .

• Consider 
$$d = 2$$
 and  $m = 3$ .  
 $k\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}\right) = (x_1y_1 + x_2y_2)^3.$ 

Feature map

$$egin{aligned} \phi\left(\left(egin{array}{c} x_1\ x_2 \end{array}
ight)
ight) &= \left(egin{array}{c} x_1^3\ \sqrt{3}x_1^2x_2\ \sqrt{3}x_1x_2^2\ x_2^3 \end{array}
ight) = \left(egin{array}{c} \phi_1(x)\ \phi_2(x)\ \phi_2(x)\ \phi_3(x)\ \phi_4(x) \end{array}
ight) \end{aligned}$$

So,  $\mathcal{H} = \mathbb{R}^4$ .

Show that  $k(x,y) = \langle \phi(x), \phi(y) 
angle_{\mathcal{H}} = \phi(x)^{ op} \phi(y).$ 

Let  $\mathcal{X} = \mathbb{R}^d$  (domain).

$$k(oldsymbol{x},oldsymbol{y}) = \left(oldsymbol{x}^{ op}oldsymbol{y}
ight)^m$$

is positive definite, for  $m \in \{1, 2, \ldots\}$ .

• Consider 
$$d = 2$$
 and  $m = 3$ .  
 $k\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}\right) = (x_1y_1 + x_2y_2)^3.$ 

Feature map

$$\phi\left(\left(egin{array}{c} x_1\ x_2\end{array}
ight)
ight)=\left(egin{array}{c} x_1^3\ \sqrt{3}x_1^2x_2\ \sqrt{3}x_1x_2^2\ x_2^3\end{array}
ight)=\left(egin{array}{c} \phi_1(x)\ \phi_2(x)\ \phi_2(x)\ \phi_3(x)\ \phi_4(x)\end{array}
ight)$$

So,  $\mathcal{H} = \mathbb{R}^4$ . Show that  $k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \phi(x)^\top \phi(y)$ .

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$egin{aligned} &k_1(m{x},m{y})+k_2(m{x},m{y})\ &k_1(f(m{x}),f(m{x}))\ &k_1(m{x},m{y})k_2(m{x},m{y})\ &\exp(k_1(m{x},m{y}))\ &dots\ &dots\$	$egin{aligned} \phi(m{x}) &=  ext{stack of } \phi_1(m{x})  ext{ and } \phi_2(m{x}) \ \phi(m{x}) &= \phi_1(f(m{x})) \ \phi(m{x}) &=  ext{tensor product of } \phi_1(m{x}), \phi_2(m{x}) \ &pprox  ext{ weighted polynomial features of all orders} \ &\vdots \end{aligned}$

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$k_1(x,y)+k_2(x,y)$	$\phi(x)= ext{stack}  ext{ of } \phi_1(x)  ext{ and } \phi_2(x)$
$k_1(f(\boldsymbol{x}),f(\boldsymbol{x}))$	$\phi(x)=\phi_1(f(x))$
$k_1({oldsymbol x},{oldsymbol y})k_2({oldsymbol x},{oldsymbol y})$	$\phi(x) =  ext{tensor}  ext{ product of } \phi_1(x), \phi_2(x)$
$\exp(k_1(x, \boldsymbol{y}))$	pprox weighted polynomial features of all orders

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$egin{aligned} &k_1(oldsymbol{x},oldsymbol{y})+k_2(oldsymbol{x},oldsymbol{y})\ &k_1(f(oldsymbol{x}),f(oldsymbol{x})) \end{aligned}$	$egin{aligned} \phi(x) &=  ext{stack of } \phi_1(x)  ext{ and } \phi_2(x) \ \phi(x) &= \phi_1(f(x)) \end{aligned}$
$egin{aligned} k_1(oldsymbol{x},oldsymbol{y})k_2(oldsymbol{x},oldsymbol{y})\ \exp(k_1(oldsymbol{x},oldsymbol{y})) \end{aligned}$	$egin{aligned} \phi(m{x}) &=  ext{tensor product of } \phi_1(m{x}), \phi_2(m{x}) \ &pprox  ext{ weighted polynomial features of all orders} \end{aligned}$

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$egin{aligned} &k_1(x,y)+k_2(x,y)\ &k_1(f(x),f(x)) \end{aligned}$	$\phi(x) =  ext{stack of } \phi_1(x)  ext{ and } \phi_2(x) \ \phi(x) = \phi_1(f(x))$
$k_1(x, y)k_2(x, y)$ $exp(k_1(x, y))$	$\phi(x) = \text{tensor product of } \phi_1(x), \phi_2(x)$
· · · · · · · · · · · · · · · · · · ·	

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$k_1(x, y) + k_2(x, y)$	$\phi(x) =  ext{stack of } \phi_1(x)  ext{ and } \phi_2(x)$
$k_1(f(x),f(x)) \ k_1(x,y)k_2(x,y)$	$arphi(x) = arphi_1(f(x)) \ \phi(x) =  ext{tensor product of } \phi_1(x), \phi_2(x)$
$\exp(k_1(x,y))$	$\approx$ weighted polynomial features of all orders
:	:

Assume k<sub>1</sub>, k<sub>2</sub> are pos. def. kernels with feature maps φ<sub>1</sub> and φ<sub>2</sub>.
New kernel k with feature map φ.

k(x,y)	features
$k_1(x, y) + k_2(x, y)$ $k_1(f(x), f(x))$	$\phi(x) =  ext{stack of } \phi_1(x)  ext{ and } \phi_2(x)$ $\phi(x) = \phi_1(f(x))$
$k_1(f(x),f(x)) \ k_1(x,y)k_2(x,y)$	$arphi(x) = arphi_1(f(x)) \ \phi(x) =  ext{tensor product of } \phi_1(x), \phi_2(x)$
$\exp(\mathit{k}_1(x,y))$	$\approx$ weighted polynomial features of all orders
:	:

# Non-Injective Mean Embedding

Variance difference revisited ....



We used  $\phi(x) = x$ . So, k(x, y) = xy (linear kernel).  $MMD_k^2(p, q) = (\mathbb{E}_{X \sim p}[\phi(X)] - \mathbb{E}_{Y \sim q}[\phi(X)])^2 = 0 \text{ but } p \neq q.$ Why?

- k (and thus  $\phi$ ) is not powerful enough.
- Mathematically, the map  $p \mapsto \mu_p$  is not injective.

# Non-Injective Mean Embedding

Variance difference revisited ....



We used  $\phi(x) = x$ . So, k(x, y) = xy (linear kernel).  $MMD_k^2(p, q) = (\mathbb{E}_{X \sim p}[\phi(X)] - \mathbb{E}_{Y \sim q}[\phi(X)])^2 = 0 \text{ but } p \neq q.$ Why?

- k (and thus  $\phi$ ) is not powerful enough.
- Mathematically, the map  $p \mapsto \mu_p$  is not injective.

# Non-Injective Mean Embedding

Variance difference revisited ....



• We used  $\phi(x) = x$ . So, k(x, y) = xy (linear kernel).

 $MMD_k^2(p,q) = \left(\mathbb{E}_{X \sim p}[\phi(X)] - \mathbb{E}_{Y \sim q}[\phi(X)]\right)^2 = 0 \text{ but } p \neq q.$ Why?

- k (and thus  $\phi$ ) is not powerful enough.
- Mathematically, the map  $p \mapsto \mu_p$  is not injective.
# Characteristic Kernels [Fukumizu et al., 2008]

### Definition 3.

A pos. def. kernel k is said to be <u>characteristic</u> if distinct distributions are embedded to different points in  $\mathcal{H}$ .

• Mathematically,  $p \mapsto \mathbb{E}_{x \sim p}[\phi(x)]$  is injective.



#### not characteristic

characteristic

- If k is characteristic, . . .
- $\mu_p$  contains all information of p,
- MMD<sub>k</sub>(p, q) = 0 if and only if p = q [Gretton et al., 2012a]. A proper distance.

# Characteristic Kernels [Fukumizu et al., 2008]

### Definition 3.

A pos. def. kernel k is said to be <u>characteristic</u> if distinct distributions are embedded to different points in  $\mathcal{H}$ .

• Mathematically,  $p \mapsto \mathbb{E}_{x \sim p}[\phi(x)]$  is injective.



not characteristic

characteristic

- If k is characteristic, ...
  - $\mu_p$  contains all information of p,
- MMD<sub>k</sub>(p, q) = 0 if and only if p = q [Gretton et al., 2012a]. A proper distance.

# Examples of Characteristic Kernels

Characteristic kernels on  $\mathcal{X} = \mathbb{R}^d$ :

Gaussian kernel:

$$k(oldsymbol{x},oldsymbol{y}) = \exp\left(-rac{\|oldsymbol{x}-oldsymbol{y}\|^2}{2\sigma^2}
ight)$$

for  $\sigma > 0$ .

• Laplace kernel:  $k(x, y) = \exp\left(-rac{\|x-y\|}{2\sigma}
ight)$  for  $\sigma > 0$ .

Matérn class of kernels [Rasmussen and Williams, 2006, Sec 4.2.1]

etc. See [Sriperumbudur et al., 2010].

#### Not characteristic:

Polynomial kernel:  $k(x, y) = (x^\top y + c)^d$  for  $c \ge 0, d \in \{1, 2, \ldots\}.$ 

## Examples of Characteristic Kernels

Characteristic kernels on  $\mathcal{X} = \mathbb{R}^d$ :

Gaussian kernel:

$$k(oldsymbol{x},oldsymbol{y}) = \exp\left(-rac{\|oldsymbol{x}-oldsymbol{y}\|^2}{2\sigma^2}
ight)$$

for  $\sigma > 0$ .

• Laplace kernel:  $k(x, y) = \exp\left(-\frac{||x-y||}{2\sigma}\right)$  for  $\sigma > 0$ .

Matérn class of kernels [Rasmussen and Williams, 2006, Sec 4.2.1]

etc. See [Sriperumbudur et al., 2010].

Not characteristic:

Polynomial kernel:  $k(x, y) = (x^\top y + c)^d$  for  $c \ge 0, d \in \{1, 2, \ldots\}.$ 

# Examples of Characteristic Kernels

Characteristic kernels on  $\mathcal{X} = \mathbb{R}^d$ :

Gaussian kernel:

$$k(oldsymbol{x},oldsymbol{y}) = \exp\left(-rac{\|oldsymbol{x}-oldsymbol{y}\|^2}{2\sigma^2}
ight)$$

for  $\sigma > 0$ .

• Laplace kernel:  $k(x, y) = \exp\left(-\frac{||x-y||}{2\sigma}\right)$  for  $\sigma > 0$ .

Matérn class of kernels [Rasmussen and Williams, 2006, Sec 4.2.1]

etc. See [Sriperumbudur et al., 2010].

#### Not characteristic:

Polynomial kernel:  $k(x, y) = (x^\top y + c)^d$  for  $c \ge 0, d \in \{1, 2, \ldots\}.$ 

Only population quantities.

- **1** Cannot compute  $\infty$ -dimensional  $\phi(\cdot)$ . Can still compute  $MMD_k(p, q)$ . Kernel trick.
- 2 Positive definite  $k(\cdot, \cdot) \iff \phi(\cdot)$  exists.
- 3 If k is characteristic,  $\mathbb{E}_{x \sim p}[\phi(x)]$  fully characterizes p.

4 Characteristic k implies  $\mathrm{MMD}_k(p,q) = \|\mathbb{E}_{{m x}\sim p}[\phi({m x})] - \mathbb{E}_{{m y}\sim q}[\phi({m y})]\| ext{ iff } p = q$ 

Only population quantities.

- 1 Cannot compute  $\infty$ -dimensional  $\phi(\cdot)$ . Can still compute  $MMD_k(p, q)$ . Kernel trick.
- 2 Positive definite  $k(\cdot, \cdot) \iff \phi(\cdot)$  exists.
- 3 If k is characteristic,  $\mathbb{E}_{x \sim p}[\phi(x)]$  fully characterizes p.



4 Characteristic k implies  $\mathrm{MMD}_k(p,q) = \|\mathbb{E}_{{m x}\sim p}[\phi({m x})] - \mathbb{E}_{{m y}\sim q}[\phi({m y})]\| ext{ iff } p = q.$ 

Only population quantities.

- 1 Cannot compute  $\infty$ -dimensional  $\phi(\cdot)$ . Can still compute  $MMD_k(p, q)$ . Kernel trick.
- 2 Positive definite  $k(\cdot, \cdot) \iff \phi(\cdot)$  exists.
- 3 If k is characteristic,  $\mathbb{E}_{x \sim p}[\phi(x)]$  fully characterizes p.



4 Characteristic k implies  $MMD_k(p,q) = \|\mathbb{E}_{\boldsymbol{x} \sim p}[\boldsymbol{\phi}(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{y} \sim q}[\boldsymbol{\phi}(\boldsymbol{y})]\|$  iff p = q.

Only population quantities.

- 1 Cannot compute  $\infty$ -dimensional  $\phi(\cdot)$ . Can still compute  $MMD_k(p, q)$ . Kernel trick.
- 2 Positive definite  $k(\cdot, \cdot) \iff \phi(\cdot)$  exists.
- 3 If k is characteristic,  $\mathbb{E}_{x \sim p}[\phi(x)]$  fully characterizes p.



4 Characteristic k implies  $MMD_k(p,q) = \|\mathbb{E}_{\boldsymbol{x} \sim p}[\boldsymbol{\phi}(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{y} \sim q}[\boldsymbol{\phi}(\boldsymbol{y})]\|$  iff p = q.



#### 1 Background

#### 2 Kernel Methods for Comparing Distributions

#### 3 Nonparametric Two-Sample Testing

#### 4 Further Topics and Conclusion

Have: Two collections of samples X, Y from unknown p and q.

 $\textbf{Goal: Test } H_0 \colon p = q \text{ vs } H_1 \colon p \neq q.$ 

1 When p = q,  $n \text{MMD}_k^2$  (random) is "close to 0."

2 When  $p \neq q$ ,  $n \text{MMD}_k^2$  is "far from 0."

Have: Two collections of samples X, Y from unknown p and q. Goal: Test  $H_0$ : p = q vs  $H_1$ :  $p \neq q$ .

1 When p = q,  $n \text{MMD}_k^2$  (random) is "close to 0."

2 When  $p \neq q$ ,  $n \text{MMD}_k^2$  is "far from 0."

Have: Two collections of samples X, Y from unknown p and q. Goal: Test  $H_0$ : p = q vs  $H_1$ :  $p \neq q$ .

- 1 When p = q,  $nMMD_k^2$  (random) is "close to 0."
- 2 When  $p \neq q$ ,  $nMMD_k^2$  is "far from 0."

Have: Two collections of samples X, Y from unknown p and q. Goal: Test  $H_0: p = q$  vs  $H_1: p \neq q$ . 1 When p = q,  $n \widehat{\text{MMD}}_k^2$  (random) is "close to 0." 2 When  $p \neq q$ ,  $n \widehat{\text{MMD}}_k^2$  is "far from 0."



Have: Two collections of samples X, Y from unknown p and q. Goal: Test  $H_0: p = q$  vs  $H_1: p \neq q$ .

- 1 When p = q,  $nMMD_k^2$  (random) is "close to 0."
- 2 When  $p \neq q$ ,  $n \text{MMD}_k^2$  is "far from 0."



# Asymptotic Null Distribution of $nMMD_k^2$

When  $H_0: p = q$ , statistic has asymptotic distribution

$$\widehat{n\mathrm{MMD}_k^2}\sim\sum_{l=1}^\infty\lambda_l\left[Z_l^2-2
ight]$$



# Asymptotic Distribution Under H<sub>1</sub> [Gretton et al., 2012a]

• When  $H_1: p \neq q$ , statistic is asymptotically normal,

$$\sqrt{n}\left(\widehat{\mathrm{MMD}_k^2}-\mathrm{MMD}_k^2(p,q)
ight) \stackrel{d}{\longrightarrow} \mathcal{N}\left(0,\,V_k(p,q)
ight),$$

 $V_k(p, q) =$ variance term.



# Gaussian kernel: $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$ . Best $\sigma^2$ ?

Keep false rejection rate at α. Maximize true rejection rate.
Keep P(reject H<sub>0</sub>|H<sub>0</sub> true) < α. Maximize P(reject H<sub>0</sub>|H<sub>1</sub> true).

Gaussian kernel: 
$$k(x, y) = \exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
. Best  $\sigma^2$ ?

• Keep false rejection rate at  $\alpha$ . Maximize true rejection rate.

• Keep  $\mathbb{P}(\text{reject } H_0|H_0 \text{ true}) \leq \alpha$ . Maximize  $\mathbb{P}(\text{reject } H_0|H_1 \text{ true})$ .

Gaussian kernel: 
$$k(x, y) = \exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
. Best  $\sigma^2$ ?

• Keep false rejection rate at  $\alpha$ . Maximize true rejection rate.

• Keep  $\mathbb{P}(\text{reject } H_0|H_0 \text{ true}) \leq \alpha$ . Maximize  $\mathbb{P}(\text{reject } H_0|H_1 \text{ true})$ .



Gaussian kernel: 
$$k(x, y) = \exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
. Best  $\sigma^2$ ?

• Keep false rejection rate at  $\alpha$ . Maximize true rejection rate.

• Keep  $\mathbb{P}(\text{reject } H_0|H_0 \text{ true}) \leq \alpha$ . Maximize  $\mathbb{P}(\text{reject } H_0|H_1 \text{ true})$ .



Gaussian kernel: 
$$k(x, y) = \exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
. Best  $\sigma^2$ ?

• Keep false rejection rate at  $\alpha$ . Maximize true rejection rate.

• Keep  $\mathbb{P}(\text{reject } H_0|H_0 \text{ true}) \leq \alpha$ . Maximize  $\mathbb{P}(\text{reject } H_0|H_1 \text{ true})$ .



Gaussian kernel: 
$$k(x, y) = \exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
. Best  $\sigma^2$ ?

• Keep false rejection rate at  $\alpha$ . Maximize true rejection rate.

• Keep  $\mathbb{P}(\text{reject } H_0|H_0 \text{ true}) \leq \alpha$ . Maximize  $\mathbb{P}(\text{reject } H_0|H_1 \text{ true})$ .



MMD Power Criterion [Sutherland et al., 2016]

• The test power  $\mathbb{P}(\text{reject } H_0 \mid H_1 \text{ true}) =$ 

$$\mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2}>\hat{c}_{lpha}
ight)
ightarrow \Phi\left(\sqrt{n}rac{\mathrm{MMD}_k^2}{\sqrt{V_k}}-rac{\hat{c}_{lpha}}{\sqrt{NV_k}}
ight),$$

where

- $\Phi$  is the CDF of the standard normal distribution.
- $\hat{c}_{\alpha}$  is an estimate of the  $1 \alpha$  quantile  $c_{\alpha}$  of the null distribution.

Choose the kernel which maximizes

 $\frac{\mathrm{MMD}_k^2}{\sqrt{V_k}}.$ 

Signal-to-noise ratio. Can be estimated with samples.

MMD Power Criterion [Sutherland et al., 2016]

• The test power  $\mathbb{P}(\text{reject } H_0 \mid H_1 \text{ true}) =$ 

$$\mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2}>\hat{c}_{lpha}
ight)
ightarrow \Phi\left(\sqrt{n}rac{\mathrm{MMD}_k^2}{\sqrt{V_k}}-rac{\hat{c}_{lpha}}{\sqrt{NV_k}}
ight),$$

where

- $\Phi$  is the CDF of the standard normal distribution.
- $\hat{c}_{\alpha}$  is an estimate of the  $1 \alpha$  quantile  $c_{\alpha}$  of the null distribution.

Choose the kernel which maximizes

 $\frac{\mathrm{MMD}_k^2}{\sqrt{V_k}}.$ 

Signal-to-noise ratio. Can be estimated with samples.

• Let  $Z \sim \mathcal{N}(0, 1)$ .

$$\begin{split} \mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2} > \widehat{c}_{\alpha}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\widehat{n\mathrm{MMD}_k^2} - n\mathrm{MMD}_k^2}{\sqrt{V_k}} > \frac{\widehat{c}_{\alpha}}{\sqrt{V_k}} - \frac{n\mathrm{MMD}_k^2}{\sqrt{V_k}}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\sqrt{n\mathrm{MMD}_k^2} - \sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} > \frac{\widehat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \mathbb{P}_{H_1}\left(Z > \frac{\widehat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \Phi\left(\frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} - \frac{\widehat{c}_{\alpha}}{\sqrt{nV_k}}\right) \end{split}$$

where

•  $\Phi$  is the CDF of the standard normal distribution.

• Let  $Z \sim \mathcal{N}(0, 1)$ .

$$\begin{split} \mathbb{P}_{H_{1}}\left(\widehat{n\mathrm{MMD}_{k}^{2}} > \widehat{c}_{\alpha}\right) \\ &= \mathbb{P}_{H_{1}}\left(\frac{\widehat{n\mathrm{MMD}_{k}^{2}} - n\mathrm{MMD}_{k}^{2}}{\sqrt{V_{k}}} > \frac{\widehat{c}_{\alpha}}{\sqrt{V_{k}}} - \frac{n\mathrm{MMD}_{k}^{2}}{\sqrt{V_{k}}}\right) \\ &= \mathbb{P}_{H_{1}}\left(\frac{\sqrt{n\mathrm{MMD}_{k}^{2}} - \sqrt{n\mathrm{MMD}_{k}^{2}}}{\sqrt{V_{k}}} > \frac{\widehat{c}_{\alpha}}{\sqrt{nV_{k}}} - \frac{\sqrt{n\mathrm{MMD}_{k}^{2}}}{\sqrt{V_{k}}}\right) \\ &\to \mathbb{P}_{H_{1}}\left(Z > \frac{\widehat{c}_{\alpha}}{\sqrt{nV_{k}}} - \frac{\sqrt{n\mathrm{MMD}_{k}^{2}}}{\sqrt{V_{k}}}\right) \\ &\to \Phi\left(\frac{\sqrt{n\mathrm{MMD}_{k}^{2}}}{\sqrt{V_{k}}} - \frac{\widehat{c}_{\alpha}}{\sqrt{nV_{k}}}\right) \end{split}$$

where

•  $\Phi$  is the CDF of the standard normal distribution.

• Let  $Z \sim \mathcal{N}(0, 1)$ .

$$\begin{split} \mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2} > \hat{c}_{\alpha}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\widehat{n\mathrm{MMD}_k^2} - n\mathrm{MMD}_k^2}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{V_k}} - \frac{n\mathrm{MMD}_k^2}{\sqrt{V_k}}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\sqrt{n\mathrm{MMD}_k^2} - \sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \mathbb{P}_{H_1}\left(Z > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \Phi\left(\frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} - \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}}\right) \end{split}$$

where

•  $\Phi$  is the CDF of the standard normal distribution.

• Let  $Z \sim \mathcal{N}(0, 1)$ .

$$\begin{split} \mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2} > \hat{c}_{\alpha}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\widehat{n\mathrm{MMD}_k^2} - n\mathrm{MMD}_k^2}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{V_k}} - \frac{n\mathrm{MMD}_k^2}{\sqrt{V_k}}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\sqrt{n\mathrm{MMD}_k^2} - \sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \mathbb{P}_{H_1}\left(Z > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \Phi\left(\frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} - \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}}\right) \end{split}$$

where

•  $\Phi$  is the CDF of the standard normal distribution.

• Let  $Z \sim \mathcal{N}(0, 1)$ .

$$\begin{split} & \mathbb{P}_{H_1}\left(\widehat{n\mathrm{MMD}_k^2} > \hat{c}_{\alpha}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\widehat{n\mathrm{MMD}_k^2} - n\mathrm{MMD}_k^2}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{V_k}} - \frac{n\mathrm{MMD}_k^2}{\sqrt{V_k}}\right) \\ &= \mathbb{P}_{H_1}\left(\frac{\sqrt{n\mathrm{MMD}_k^2} - \sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \mathbb{P}_{H_1}\left(Z > \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}} - \frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}}\right) \\ &\to \Phi\left(\frac{\sqrt{n\mathrm{MMD}_k^2}}{\sqrt{V_k}} - \frac{\hat{c}_{\alpha}}{\sqrt{nV_k}}\right) \end{split}$$

where

•  $\Phi$  is the CDF of the standard normal distribution.

•  $\hat{c}_{\alpha}$  is an estimate of the  $1 - \alpha$  quantile  $c_{\alpha}$  of the null distribution.

33/44

# Properties of the MMD Test

- As sample size  $n \to \infty$ ,
- 1 If  $H_0: p = q$ , then  $\mathbb{P}(\text{reject } H_0) \leq \alpha$ .
- 2 If k is characteristic and  $H_1 \colon p \neq q$ , then  $\mathbb{P}(\text{reject } H_0) \to 1$ .
- (1) and (2)  $\implies$  a consistent test.
- $MMD_k^2$  can be estimated in  $\mathcal{O}(n^2)$  time.
  - But, linear-time versions (O(n)) exist [Gretton et al., 2012a, Zaremba et al., 2013].

# Properties of the MMD Test

As sample size  $n \to \infty$ ,

- 1 If  $H_0: p = q$ , then  $\mathbb{P}(\text{reject } H_0) \leq \alpha$ .
- 2 If k is characteristic and  $H_1: p \neq q$ , then  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$ .
- (1) and (2)  $\implies$  a consistent test.
- $MMD_k^2$  can be estimated in  $\mathcal{O}(n^2)$  time.
  - But, linear-time versions (O(n)) exist [Gretton et al., 2012a, Zaremba et al., 2013].

# Properties of the MMD Test

As sample size  $n \to \infty$ ,

- 1 If  $H_0: p = q$ , then  $\mathbb{P}(\text{reject } H_0) \leq \alpha$ .
- 2 If k is characteristic and  $H_1: p \neq q$ , then  $\mathbb{P}(\text{reject } H_0) \rightarrow 1$ .
- (1) and (2)  $\implies$  a consistent test.
- $MMD_k^2$  can be estimated in  $\mathcal{O}(n^2)$  time.
  - But, linear-time versions (O(n)) exist [Gretton et al., 2012a, Zaremba et al., 2013].

# Generate MNIST Handwritten Digits





Observed MNIST handwritten digits. X.

Generated images from a model.  $\checkmark$ .

• Goal: Learn a function which transforms noise into a handwritten digit.

# Generative Moment Matching Networks [Li et al., 2015]

#### **Generative Moment Matching Networks**

Yujia Li<sup>1</sup> YUJIALI@CS.TORONTO.EDU Kevin Swersky<sup>1</sup> KSWERSKY@CS.TORONTO.EDU Richard Zemel<sup>1,2</sup> ZEMEL@CS.TORONTO.EDU <sup>1</sup>Department of Computer Science, University of Toronto, ON, CANADA

<sup>2</sup>Canadian Institute for Advanced Research, Toronto, ON, CANADA

- ICML 2015.
- Code: https://github.com/yujiali/gmmn
- One of the first to use MMD to train a generative network.

# More Recent Works on MMD Based Generative Nets

MMD GAN: Towards Deeper Understanding of Moment Matching Network Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, Barnabás Póczos https://arxiv.org/abs/1705.08584

Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy Dougal J. Sutherland, Hsiao-Yu Tung, Heiko Strathmann, Soumyajit De, Aaditya Ramdas, Alex Smola, Arthur Gretton ICLR 2017 https://arxiv.org/abs/1611.04488

Demystifying MMD GANs Mikolaj Binkowski, Dougal J. Sutherland, Michael Arbel, Arthur Gretton ICLR 2018

https://openreview.net/pdf?id=r11UOzWCW
Generative Moment Matching Networks [Li et al., 2015]

 $\arg\min_{\boldsymbol{\alpha}} \operatorname{MMD}_{k}^{2}(\mathsf{X}, \{g_{\theta}(\boldsymbol{z}_{i})\}_{i=1}^{n})$ 

- X = training sets. x<sub>i</sub> = one digit (an image with 28 × 28 = 784 pixels). 60000 images.
- **Z** =  $\{z_i\}_{i=1}^n$  random noise vectors. Drawn from  $\mathcal{N}(0, I)$ .
- $g_{\theta}(z)$  a deep net transforming noise z into an image.
- Kernel k: sum of 5 Gaussian kernels of different bandwidths

Network architecture (my own, not [Li et al., 2015]):

- 4 hidden layers. Total parameters 60,608 (in  $\theta$ ).
- Training for 15 epochs.  $\approx$  7 minutes. My laptop without GPU.

Generative Moment Matching Networks [Li et al., 2015]

 $\arg\min_{\boldsymbol{\alpha}} \operatorname{MMD}_{k}^{2}(\mathsf{X}, \{g_{\theta}(\boldsymbol{z}_{i})\}_{i=1}^{n})$ 

- X = training sets. x<sub>i</sub> = one digit (an image with 28 × 28 = 784 pixels). 60000 images.
- $Z = \{z_i\}_{i=1}^n$  random noise vectors. Drawn from  $\mathcal{N}(0, I)$ .
- $g_{\theta}(z)$  a deep net transforming noise z into an image.
- Kernel k: sum of 5 Gaussian kernels of different bandwidths

Network architecture (my own, not [Li et al., 2015]):

- 4 hidden layers. Total parameters 60,608 (in  $\theta$ ).
- Training for 15 epochs.  $\approx$  7 minutes. My laptop without GPU.

Generative Moment Matching Networks [Li et al., 2015]

 $\arg\min_{\boldsymbol{\alpha}} \operatorname{MMD}_{k}^{2}(\mathsf{X}, \{g_{\theta}(\boldsymbol{z}_{i})\}_{i=1}^{n})$ 

- X = training sets. x<sub>i</sub> = one digit (an image with 28 × 28 = 784 pixels). 60000 images.
- **Z** =  $\{z_i\}_{i=1}^n$  random noise vectors. Drawn from  $\mathcal{N}(0, I)$ .
- $g_{\theta}(z)$  a deep net transforming noise z into an image.
- Kernel k: sum of 5 Gaussian kernels of different bandwidths

Network architecture (my own, not [Li et al., 2015]):

- 4 hidden layers. Total parameters 60,608 (in  $\theta$ ).
- Training for 15 epochs.  $\approx$  7 minutes. My laptop without GPU.

#### My Results



 $\mathbf{39}/\mathbf{44}$ 

# Quick Comments



(a) GMMN MNIST samples



(b) GMMN TFD samples





(c) GMMN+AE MNIST samples (d) GMMN+AE TFD samples





(f) GMMN+AE nearest neighbors for MNIST samples



(g) GMMN nearest neighbors for TFD samples

101	200	23	25	1	5	ele:	E.
15.4		22	35	101	25		154

(h) GMMN+AE nearest neighbors for TFD samples

- I could have done better. Just had to wait + bigger network. Key points:
  - Easy to train. Simple implementation.
  - Stable training.
  - Image quality depends on kernel k.



#### 1 Background

2 Kernel Methods for Comparing Distributions

3 Nonparametric Two-Sample Testing

4 Further Topics and Conclusion

# Further Topics I

"Dual view": Reproducing Kernel Hilbert Spaces (RKHSs)

• Each point in  $\mathcal{H}$  can be seen as a function:

 $f\in \mathcal{H}\iff f(x)=\sum_{i=1}^m lpha_i k(x,x_i) ext{ for some } \{lpha_i\}_{i=1}^m, \{x_i\}_{i=1}^m.$ MMD

- Associated with MMD(p, q) is the witness function.
- Unit-norm function in  $\mathcal{H}$  that best distinguishes p and q.

### Further Topics I

"Dual view": Reproducing Kernel Hilbert Spaces (RKHSs)

Each point in  $\mathcal{H}$  can be seen as a function:

 $f\in \mathcal{H}\iff f(x)=\sum_{i=1}^m lpha_i k(x,x_i) ext{ for some } \{lpha_i\}_{i=1}^m, \{x_i\}_{i=1}^m.$ MMD



Associated with MMD(p, q) is the witness function.
Unit-norm function in H that best distinguishes p and q.

# Further Topics II

#### Dependence measure

- Recall X independent of Y iff  $p_{xy}(X, Y) = p_x(X)p_y(Y)$ .
- MMD(*p<sub>xy</sub>*, *p<sub>x</sub>p<sub>y</sub>*) can be used to measure dependence [Gretton et al., 2005].
- <u>Applications</u>: Feature selection, clustering etc.

Others

- Linear-time versions of MMD [Gretton et al., 2012b, Chwialkowski et al., 2015, Jitkrittum et al., 2016].
- Goodness-of-fit test by distance(model, data)
   [Liu et al., 2016, Chwialkowski et al., 2016, Jitkrittum et al., 2017].
- Gaussian process regression/classification
   [Rasmussen and Williams, 2006].

# Further Topics II

#### Dependence measure

- Recall X independent of Y iff  $p_{xy}(X, Y) = p_x(X)p_y(Y)$ .
- MMD(*p<sub>xy</sub>*, *p<sub>x</sub>p<sub>y</sub>*) can be used to measure dependence [Gretton et al., 2005].
- Applications: Feature selection, clustering etc.

Others

- Linear-time versions of MMD [Gretton et al., 2012b, Chwialkowski et al., 2015, Jitkrittum et al., 2016].
- Goodness-of-fit test by distance(model, data)
   [Liu et al., 2016, Chwialkowski et al., 2016, Jitkrittum et al., 2017].
- Gaussian process regression/classification
   [Rasmussen and Williams, 2006].

# Conclusion



 Maximum Mean Discrepancy (MMD) = distance between two distributions

- "Mean embed" distributions to a high-dimensional space  $\mathcal{H}$ .
- Measure the distance in  $\mathcal{H}$ .
- Characteristic kernel (e.g., Gaussian kernel)
  - $\implies \mathrm{MMD}(p,q) = \mathsf{0} ext{ iff } p = q.$
- **Two-sample testing with MMD. Consistent test.**

# Conclusion



 Maximum Mean Discrepancy (MMD) = distance between two distributions

- "Mean embed" distributions to a high-dimensional space  $\mathcal{H}$ .
- Measure the distance in  $\mathcal{H}$ .
- Characteristic kernel (e.g., Gaussian kernel)

 $\implies$  MMD(p,q) = 0 iff p = q.

**Two-sample testing with MMD. Consistent test.** 

# Conclusion



 Maximum Mean Discrepancy (MMD) = distance between two distributions

- "Mean embed" distributions to a high-dimensional space  $\mathcal{H}$ .
- Measure the distance in  $\mathcal{H}$ .
- Characteristic kernel (e.g., Gaussian kernel)

 $\implies$  MMD(p, q) = 0 iff p = q.

**Two-sample testing** with MMD. Consistent test.



# Thank you

Wittawat Jitkrittum

wittawat.com

wittawatj@gmail.com

45/44

### References I

Chwialkowski, K., Ramdas, A., Sejdinovic, D., and Gretton, A. (2015).

Fast two-sample testing with analytic representations of probability measures.

In NIPS, pages 1972–1980.

- Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).
   A kernel test of goodness of fit.
   In *ICML*, pages 2606-2615.
- Fukumizu, K., Gretton, A., Sun, X., and Schölkopf, B. (2008). Kernel measures of conditional dependence. In NIPS, pages 489–496.

#### References II

Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., and Smola, A. (2012a).

A kernel two-sample test.

Journal of Machine Learning Research, 13:723–773.

 Gretton, A., Herbrich, R., Smola, A., Bousquet, O., and Schölkopf, B. (2005).
 Kernel methods for measuring independence.
 Journal of Machine Learning Research, 6:2075-2129.

 Gretton, A., Sejdinovic, D., Strathmann, H., Balakrishnan, S., Pontil, M., Fukumizu, K., and Sriperumbudur, B. K. (2012b).
 Optimal kernel choice for large-scale two-sample tests.
 In NIPS, pages 1205-1213.

### References III

Jitkrittum, W., Szabó, Z., Chwialkowski, K. P., and Gretton, A. (2016).

Interpretable Distribution Features with Maximum Testing Power.

In *NIPS*, pages 181–189.



Jitkrittum, W., Xu, W., Szabo, Z., Fukumizu, K., and Gretton, A. (2017).

A linear-time kernel goodness-of-fit test.

Li, Y., Swersky, K., and Zemel, R. (2015). Generative moment matching networks. In *ICML*, pages 1718–1727.

### References IV

- Liu, Q., Lee, J., and Jordan, M. (2016).
   A kernelized Stein discrepancy for goodness-of-fit tests.
   In *ICML*, pages 276-284.
- Rasmussen, C. E. and Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. MIT Press, Cambridge, MA.
- Shawe-Taylor, J. and Cristianini, N. (2004). Kernel methods for pattern analysis. Cambridge university press.
- Smola, A., Gretton, A., Song, L., and Schölkopf, B. (2007).
   A Hilbert space embedding for distributions.
   In International Conference on Algorithmic Learning Theory (ALT), pages 13-31.

#### References V

- Sriperumbudur, B., Gretton, A., Fukumizu, K., Schoelkopf, B., and Lanckriet, G. (2010).
   Hilbert space embeddings and metrics on probability measures. Journal of Machine Learning Research, 11:1517-1561.
- Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ramdas, A., Smola, A., and Gretton, A. (2016).
   Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy. arXiv: 1611.04488.
- Zaremba, W., Gretton, A., and Blaschko, M. (2013).
   B-test: A non-parametric, low variance kernel two-sample test.
   In NIPS, pages 755-763.