

Kernel Machines

Alain Celisse

SAMM

Paris 1-Panthéon Sorbonne University

`alain.celisse@univ-paris1.fr`

Lecture 5: Maximum Mean Discrepancy

Master 2 Data Science – Centrale Lille, Lille University
Fall 2022

Successive topics of the coming lectures:

1. Introduction to Kernel methods
2. Support vector classifiers and Kernel methods
3. Extending classical strategies to high dimension
 - ▶ KRR/LS-SVMs
 - ▶ KPCA
4. Duality gap and KKT conditions
5. Designing reproducing kernels
6. Maximum Mean Discrepancy (MMD) (Today!)

Outline of the lecture

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

- ▶ Distance between Probability distributions
- ▶ Maximum Mean Discrepancy (MMD)
- ▶ Mean embedding
- ▶ Estimation
- ▶ Two-sample test

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Distance between Probability distributions

Example: Two-sample test

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

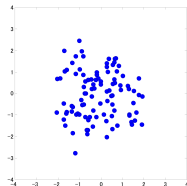
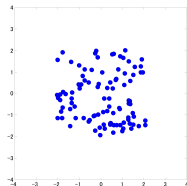
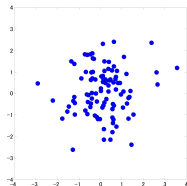
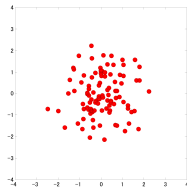
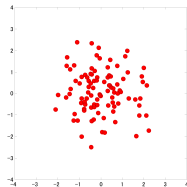
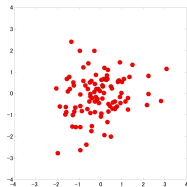
Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test



Example: Two-sample test

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

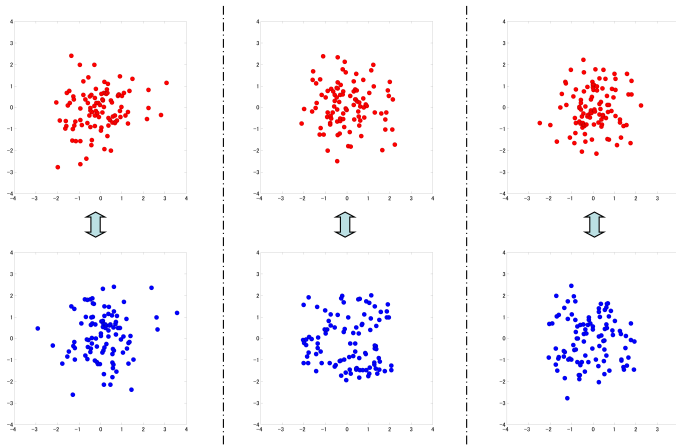
Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test



Question:

Are red and blue samples drawn from the same distribution?

Example: Two-sample test

Two-sample test

$X_1, \dots, X_n \sim P_X$ and $Y_1, \dots, Y_m \sim P_Y$, with $P_X = P_Y$?

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Example: Two-sample test

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Two-sample test

$X_1, \dots, X_n \sim P_X$ and $Y_1, \dots, Y_m \sim P_Y$, with $P_X = P_Y$?

Connections:

- Independence testing
(comparison between $P_X \otimes P_X$ and $P_{X,Y}$)

Example: Two-sample test

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Two-sample test

$X_1, \dots, X_n \sim P_X$ and $Y_1, \dots, Y_m \sim P_Y$, with $P_X = P_Y$?

Connections:

- ▶ Independence testing
(comparison between $P_X \otimes P_X$ and $P_{X,Y}$)
- ▶ Novelty/Anomaly detection (on-line)

Example: Two-sample test

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Two-sample test

$X_1, \dots, X_n \sim P_X$ and $Y_1, \dots, Y_m \sim P_Y$, with $P_X = P_Y$?

Connections:

- ▶ Independence testing
(comparison between $P_X \otimes P_X$ and $P_{X,Y}$)
- ▶ Novelty/Anomaly detection (on-line)
- ▶ Change-point detection (see the next slides!)

Example: Two-sample test

Two-sample test

$X_1, \dots, X_n \sim P_X$ and $Y_1, \dots, Y_m \sim P_Y$, with $P_X = P_Y$?

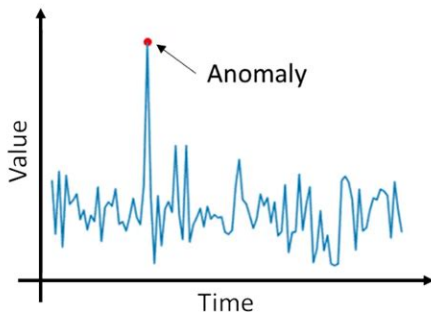
Connections:

- ▶ Independence testing
(comparison between $P_X \otimes P_X$ and $P_{X,Y}$)
- ▶ Novelty/Anomaly detection (on-line)
- ▶ Change-point detection (see the next slides!)

Key ingredient

Requires a “distance” between P_X and P_Y

Anomaly detection: Online scenario



- ▶ Reference distribution until time $t_0 > 0$
- ▶ **Goal:** From $t > t_0$, detecting (in real-time) potential shifts in the distribution

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

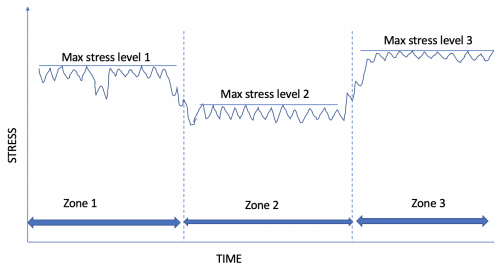
Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Change-point detection: Offline scenario



Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Change-point detection: Offline scenario

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

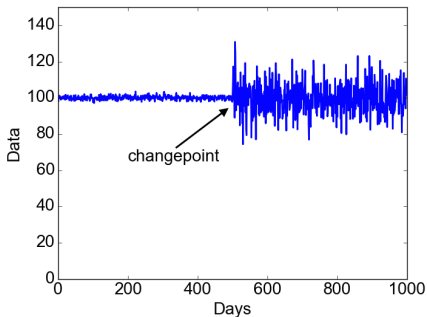
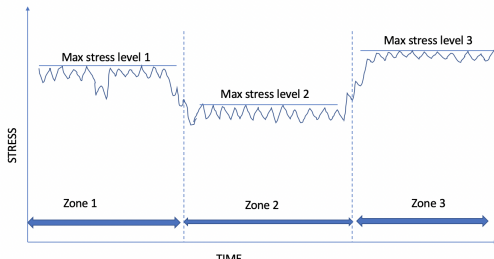
Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test



Distance between Probability distributions

Kolmogorov-Smirnov distance (KS-distance)

- ▶ X_1, \dots, X_n , and $Y_1, \dots, Y_m \in \mathbb{R}$
- ▶ Cumulative distribution functions are **characteristic**:

Theorem

$$F(t) = \mathbb{P}[X \leq t] = G(t) = \mathbb{P}[Y \leq t], \quad \forall t \in \mathbb{R}$$

implies that $P_X = P_Y$

KS-distance between P and Q

Definition (KS-distance)

$$\begin{aligned} d_{KS}(P, Q) &= \sup_{t \in \mathbb{R}} |F(t) - G(t)| = \|F - G\|_{\infty} \\ &= \sup_{t \in \mathbb{R}} |\mathbb{E}[\mathbb{1}_{(X \leq t)}] - \mathbb{E}[\mathbb{1}_{(Y \leq t)}]| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]| \\ \mathcal{F} &= \{f = \mathbb{1}_{]-\infty, t]} \mid t \in \mathbb{R}\} \end{aligned}$$

Other famous metrics

Integral probability metrics

For some class \mathcal{F} ,

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Other famous metrics

Integral probability metrics

For some class \mathcal{F} ,

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Total variation (TV) distance

$$\mathcal{F} = \mathcal{F}_{TV} = \{f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_{\infty} \leq 1\}$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Other famous metrics

Integral probability metrics

For some class \mathcal{F} ,

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Total variation (TV) distance

$$\mathcal{F} = \mathcal{F}_{TV} = \{f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_{\infty} \leq 1\}$$

Wasserstein (W) distance

$$\mathcal{F} = \mathcal{F}_W = \left\{ f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_L = \sup_{x \neq y} \left| \frac{f(x) - f(y)}{\rho(x, y)} \right| \leq 1 \right\}$$

Other famous metrics

Integral probability metrics

For some class \mathcal{F} ,

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Total variation (TV) distance

$$\mathcal{F} = \mathcal{F}_{TV} = \{f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_{\infty} \leq 1\}$$

Wasserstein (W) distance

$$\mathcal{F} = \mathcal{F}_W = \left\{ f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_L = \sup_{x \neq y} \left| \frac{f(x) - f(y)}{\rho(x, y)} \right| \leq 1 \right\}$$

L^q -distance

$$\mathcal{F} = \mathcal{F}_q = \left\{ f \in \mathcal{M}(\mathcal{X}) \mid \|f\|_q^q = \int_{\mathcal{X}} |f|^q d\lambda \leq 1 \right\}$$
$$(1 \leq q < +\infty)$$

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum
- ▶ \mathcal{F} is infinite

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum
- ▶ \mathcal{F} is infinite
- ▶ $d_{\mathcal{F}}(P, Q)$ usually difficult to compute

Distance between
Probability
distributions

Motivating example:
Two-sample test

Metric over probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum
- ▶ \mathcal{F} is infinite
- ▶ $d_{\mathcal{F}}(P, Q)$ usually difficult to compute

Richness of \mathcal{F}

- ▶ The richness of \mathcal{F} determines the properties of $d_{\mathcal{F}}(\cdot, \cdot)$

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum
- ▶ \mathcal{F} is infinite
- ▶ $d_{\mathcal{F}}(P, Q)$ usually difficult to compute

Richness of \mathcal{F}

- ▶ The richness of \mathcal{F} determines the properties of $d_{\mathcal{F}}(\cdot, \cdot)$
- ▶ **Ex:**
Is the next set characteristic?

$$\mathcal{F} = \left\{ \theta |\cdot|^2 \mid \theta \in \mathbb{R} \right\}$$

Computations

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

- ▶ Previous metrics defined from a supremum
- ▶ \mathcal{F} is infinite
- ▶ $d_{\mathcal{F}}(P, Q)$ usually difficult to compute

Richness of \mathcal{F}

- ▶ The richness of \mathcal{F} determines the properties of $d_{\mathcal{F}}(\cdot, \cdot)$
- ▶ **Ex:**
Is the next set characteristic?

$$\mathcal{F} = \left\{ \theta |\cdot|^2 \mid \theta \in \mathbb{R} \right\}$$

→ Infinite but too much thin for spanning a metric!

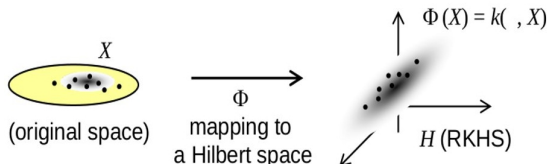
Maximum Mean Discrepancy (MMD)

Mapping data from \mathcal{X} to \mathcal{H}

- ▶ $k(\cdot, \cdot)$: psd kernel (reproducing)
- ▶ $x \in \mathbb{R}^d \mapsto k_x = k(x, \cdot) \in \mathcal{H}$: canonical feature map from \mathbb{R}^d to \mathcal{H} with

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}$$

- ▶ $k_x = \phi(x) \in \mathcal{H}$: new “observation”



Theorem (From psd to reproducing kernel)

Any psd kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ gives rise to a unique Hilbert space endowed with the scalar product

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

and such that

Theorem (From psd to reproducing kernel)

Any psd kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ gives rise to a unique Hilbert space endowed with the scalar product

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

and such that

1. \mathcal{H} contains all functions $k_x : x \mapsto k(x, \cdot)$, for all $x \in \mathcal{X}$

Theorem (From psd to reproducing kernel)

Any psd kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ gives rise to a unique Hilbert space endowed with the scalar product

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

and such that

1. \mathcal{H} contains all functions $k_x : x \mapsto k(x, \cdot)$, for all $x \in \mathcal{X}$
2. $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}} \quad (\text{Reproducing property})$$

Theorem (From psd to reproducing kernel)

Any psd kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ gives rise to a unique Hilbert space endowed with the scalar product

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

and such that

1. \mathcal{H} contains all functions $k_x : x \mapsto k(x, \cdot)$, for all $x \in \mathcal{X}$
2. $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}} \quad (\text{Reproducing property})$$

Then,

- $\mathcal{H} = \mathcal{H}_k$: Reproducing Kernel Hilbert Space (RKHS) associated with k .

Theorem (From psd to reproducing kernel)

Any psd kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ gives rise to a unique Hilbert space endowed with the scalar product

$$k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

and such that

1. \mathcal{H} contains all functions $k_x : x \mapsto k(x, \cdot)$, for all $x \in \mathcal{X}$
2. $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}} \quad (\text{Reproducing property})$$

Then,

- ▶ $\mathcal{H} = \mathcal{H}_k$: Reproducing Kernel Hilbert Space (RKHS) associated with k .
- ▶ k : Reproducing kernel of \mathcal{H} .

Maximum Mean Discrepancy (MMD)

(1/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

Definition (MMD)

(Gretton, Fukumizu, Sriperumbudur, . . .)

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Maximum Mean Discrepancy (MMD)

(1/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

Definition (MMD)

(Gretton, Fukumizu, Sriperumbudur, . . .)

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

where

$$\mathcal{F} = \{f \in \mathcal{H} \mid \|f\|_{\mathcal{H}} \leq 1\}$$

Maximum Mean Discrepancy (MMD)

(2/2)

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

Maximum Mean Discrepancy (MMD)

(2/2)

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Calculating the MMD

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

Maximum Mean Discrepancy (MMD)

(2/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Calculating the MMD

$$\begin{aligned} MMD_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[\langle f, k_X - k_Y \rangle_{\mathcal{H}}]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}_{X,Y}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}_{X,Y}[k_X - k_Y]\|_{\mathcal{H}}'' \end{aligned}$$

Maximum Mean Discrepancy (MMD)

(2/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Calculating the MMD

$$\begin{aligned} MMD_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[\langle f, k_X - k_Y \rangle_{\mathcal{H}}]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}_{X,Y}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}_{X,Y}[k_X - k_Y]\|_{\mathcal{H}}'' \end{aligned}$$

Remarks:

- No supremum anymore!

Maximum Mean Discrepancy (MMD)

(2/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Calculating the MMD

$$\begin{aligned} MMD_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[\langle f, k_X - k_Y \rangle_{\mathcal{H}}]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}_{X,Y}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}_{X,Y}[k_X - k_Y]\|''_{\mathcal{H}} \end{aligned}$$

Remarks:

- ▶ No supremum anymore!
- ▶ Meaning of $\mathbb{E}_{X,Y}[k_X - k_Y]$?

Maximum Mean Discrepancy (MMD)

(2/2)

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

$$MMD_{\mathcal{H}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]|$$

Calculating the MMD

$$\begin{aligned} MMD_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{E}_{X,Y}[\langle f, k_X - k_Y \rangle_{\mathcal{H}}]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}_{X,Y}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}_{X,Y}[k_X - k_Y]\|''_{\mathcal{H}} \end{aligned}$$

Remarks:

- ▶ No supremum anymore!
- ▶ Meaning of $\mathbb{E}_{X,Y}[k_X - k_Y]$? → Mean element

Definition (Mean element)

- ▶ k : psd kernel on \mathcal{X}
- ▶ Assume: $\mathbb{E} \left[\sqrt{k(X, X)} \right] = \mathbb{E} [\|k_X\|_{\mathcal{H}}] < +\infty$

Definition (Mean element)

- ▶ k : psd kernel on \mathcal{X}
- ▶ Assume: $\mathbb{E} \left[\sqrt{k(X, X)} \right] = \mathbb{E} [\|k_X\|_{\mathcal{H}}] < +\infty$

Then, there **exists a unique** element $\mu_P \in \mathcal{H}$ such that

$$\langle \mu_P, f \rangle_{\mathcal{H}} = \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}], \quad \forall f \in \mathcal{H}$$

Definition (Mean element)

- ▶ k : psd kernel on \mathcal{X}
- ▶ Assume: $\mathbb{E} \left[\sqrt{k(X, X)} \right] = \mathbb{E} [\|k_X\|_{\mathcal{H}}] < +\infty$

Then, there **exists a unique** element $\mu_P \in \mathcal{H}$ such that

$$\langle \mu_P, f \rangle_{\mathcal{H}} = \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}], \quad \forall f \in \mathcal{H}$$

Remark: $\mathbb{E} \left[\sqrt{k(X, X)} \right] < +\infty$ true for any bounded kernel
(e.g. normalized kernel)

Definition (Mean element)

- ▶ k : psd kernel on \mathcal{X}
- ▶ Assume: $\mathbb{E} \left[\sqrt{k(X, X)} \right] = \mathbb{E} [\|k_X\|_{\mathcal{H}}] < +\infty$

Then, there **exists a unique** element $\mu_P \in \mathcal{H}$ such that

$$\langle \mu_P, f \rangle_{\mathcal{H}} = \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}], \quad \forall f \in \mathcal{H}$$

Remark: $\mathbb{E} \left[\sqrt{k(X, X)} \right] < +\infty$ true for any bounded kernel
(e.g. normalized kernel)

Proof.

Definition (Mean element)

- ▶ k : psd kernel on \mathcal{X}
- ▶ Assume: $\mathbb{E} \left[\sqrt{k(X, X)} \right] = \mathbb{E} [\|k_X\|_{\mathcal{H}}] < +\infty$

Then, there **exists a unique** element $\mu_P \in \mathcal{H}$ such that

$$\langle \mu_P, f \rangle_{\mathcal{H}} = \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}], \quad \forall f \in \mathcal{H}$$

Remark: $\mathbb{E} \left[\sqrt{k(X, X)} \right] < +\infty$ true for any bounded kernel
(e.g. normalized kernel)

Proof.

$f \mapsto \mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}]$ is a bounded linear form over \mathcal{H} since

$$|\mathbb{E} [\langle k_X, f \rangle_{\mathcal{H}}]| \leq \mathbb{E} [\|k_X\|_{\mathcal{H}}] \|f\|_{\mathcal{H}}$$

Riesz's theorem yields the existence and unicity



Theorem

$$\blacktriangleright \mu_P = \mathbb{E}_X [k_X] \in \mathcal{H}$$

Theorem

► $\mu_P = \mathbb{E}_X [k_X] \in \mathcal{H}$

► For all $x \in \mathcal{X}$,

$$\mu_P(x) = \mathbb{E}_X [k_X(x)] = \mathbb{E}_X [k(X, x)] = \int_{\mathcal{X}} k(u, x) dP_X(u)$$

Theorem

► $\mu_P = \mathbb{E}_X [k_X] \in \mathcal{H}$

► For all $x \in \mathcal{X}$,

$$\mu_P(x) = \mathbb{E}_X [k_X(x)] = \mathbb{E}_X [k(X, x)] = \int_{\mathcal{X}} k(u, x) dP_X(u)$$

► $\mathbb{E}_X [k_X] + \mathbb{E}_Y [k_Y] = \mu_P + \mu_Q$

MMD and mean element

$$\begin{aligned} \text{MMD}_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}[k_X - k_Y]\|_{\mathcal{H}} \\ &= \|\mu_P - \mu_Q\|_{\mathcal{H}} \end{aligned}$$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

MMD expression

Mean Embedding/Mean
Element

Mean element

Two-sample test

MMD and mean element

$$\begin{aligned} \text{MMD}_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}[k_X - k_Y]\|_{\mathcal{H}} \\ &= \|\mu_P - \mu_Q\|_{\mathcal{H}} \end{aligned}$$

The “distance” between P and Q translates into a difference between the mean elements of μ_P and μ_Q in $\|\cdot\|_{\mathcal{H}}$

MMD and mean element

$$\begin{aligned} \text{MMD}_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}[k_X - k_Y]\|_{\mathcal{H}} \\ &= \|\mu_P - \mu_Q\|_{\mathcal{H}} \end{aligned}$$

The “distance” between P and Q translates into a difference between the mean elements of μ_P and μ_Q in $\|\cdot\|_{\mathcal{H}}$

Remark:

$$P = Q \quad \Rightarrow \quad \|\mu_P - \mu_Q\|_{\mathcal{H}} = 0 \quad \Leftrightarrow \quad \mu_P = \mu_Q$$

MMD and mean element

$$\begin{aligned} \text{MMD}_{\mathcal{H}}(P, Q) &= \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X) - f(Y)]| \\ &= \sup_{f \in \mathcal{F}} |\langle f, \mathbb{E}[k_X - k_Y] \rangle_{\mathcal{H}}| \\ &= \|\mathbb{E}[k_X - k_Y]\|_{\mathcal{H}} \\ &= \|\mu_P - \mu_Q\|_{\mathcal{H}} \end{aligned}$$

The “distance” between P and Q translates into a difference between the mean elements of μ_P and μ_Q in $\|\cdot\|_{\mathcal{H}}$

Remark:

$$P = Q \quad \Rightarrow \quad \|\mu_P - \mu_Q\|_{\mathcal{H}} = 0 \quad \Leftrightarrow \quad \mu_P = \mu_Q$$

Warning: The converse is not true in general!

Mean element/Mean embedding

Polynomial kernel and mean element

Mean element

For $X \sim P$

$$\begin{aligned}\mu_P &= \mathbb{E}[k_X] = \mathbb{E}[k(X, \cdot)] \in \mathcal{H} \\ \Rightarrow \mu_P(t) &= \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}\end{aligned}$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

Polynomial kernel and mean element

Kernel Machines

Alain Celisse

Mean element

For $X \sim P$

$$\begin{aligned}\mu_P &= \mathbb{E}[k_X] = \mathbb{E}[k(X, \cdot)] \in \mathcal{H} \\ \Rightarrow \mu_P(t) &= \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}\end{aligned}$$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

Examples: Polynomial kernels

- $k(X, t) = (Xt)^d$, for $X, t \in \mathbb{R}$ and $d \in \mathbb{N}^*$

$$\mu_P(t) = \mathbb{E}[(Xt)^d] = \mathbb{E}[X^d] t^d$$

→ only involves the **d th moment**.

Mean element

For $X \sim P$

$$\begin{aligned}\mu_P &= \mathbb{E}[k_X] = \mathbb{E}[k(X, \cdot)] \in \mathcal{H} \\ \Rightarrow \mu_P(t) &= \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}\end{aligned}$$

Examples: Polynomial kernels

- ▶ $k(X, t) = (Xt)^d$, for $X, t \in \mathbb{R}$ and $d \in \mathbb{N}^*$

$$\mu_P(t) = \mathbb{E}[(Xt)^d] = \mathbb{E}[X^d] t^d$$

→ only involves the **dth moment**.

- ▶ $k(X, t) = (Xt + c)^d$, for $X, t \in \mathbb{R}$ and $d \in \mathbb{N}^*$ $c > 0$

$$\mu_P(t) = \mathbb{E}[(Xt + c)^d] = \sum_{i=0}^d \binom{d}{i} \mathbb{E}[X^i] t^i c^{d-i}$$

→ involves **all moments up to d!**

Taylor expansion and mean element

For $X \sim P$

$$\mu_P(t) = \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}$$

Examples: Taylor expansion

Assume

$$k(X, t) = \sum_{i=0}^{+\infty} a_i X^i t^i, \quad \forall X, t \in \mathcal{X}$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

Taylor expansion and mean element

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

For $X \sim P$

$$\mu_P(t) = \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}$$

Examples: Taylor expansion

Assume

$$k(X, t) = \sum_{i=0}^{+\infty} a_i X^i t^i, \quad \forall X, t \in \mathcal{X}$$

Then

$$\mu_P(t) = \sum_{i=0}^{+\infty} a_i \mathbb{E}[X^i] t^i,$$

Taylor expansion and mean element

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

For $X \sim P$

$$\mu_P(t) = \mathbb{E}[k(X, t)], \quad \forall t \in \mathcal{X}$$

Examples: Taylor expansion

Assume

$$k(X, t) = \sum_{i=0}^{+\infty} a_i X^i t^i, \quad \forall X, t \in \mathcal{X}$$

Then

$$\mu_P(t) = \sum_{i=0}^{+\infty} a_i \mathbb{E}[X^i] t^i,$$

Remark:

- ▶ The exponential kernel
- ▶ All moments involved within the mean element
- ▶ Construct kernel with prescribed moments. . .

Laplace transform

$$t \in \mathcal{X} \mapsto \mathcal{L}_P(t) = \mathbb{E} \left[e^{-\langle X, t \rangle} \right]$$

- \mathcal{L}_P is characteristic of the probability distribution P

$$\mathcal{L}_P = \mathcal{L}_Q \iff P = Q$$

- Since $t \mapsto \mathcal{L}_P(t) = \mu_P(t)$ with the exponential kernel, we have

$$\mu_P = \mu_Q \Rightarrow P = Q$$

Remark:

Analogy with the characteristic function

$$t \in \mathcal{X} \mapsto \phi_P(t) = \mathbb{E} \left[e^{i\langle X, t \rangle} \right]$$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

Definition (Characteristic kernel)

A kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is characteristic if the mapping $P \mapsto \mu_P$ is injective that is,

$$\mu_P = \mu_Q \quad \Rightarrow \quad P = Q$$

Remark:

- ▶ $P \mapsto \mu_P$ justifies the name “mean embedding”
- ▶ characteristicity through mean element
- ▶ the converse inequality holds always true!
→ $P \mapsto \mu_P$: one-to-one mapping (characteristic kernel)

Proposition

The following claims are equivalent:

1. k : characteristic
2. $P \mapsto \mu_P$: injective
3. $(\mathbb{E}[f(X)] = \mathbb{E}[f(Y)], \quad \forall f \in \mathcal{H}) \Rightarrow P = Q$

Proof.

Do it!



General domain \mathcal{X}

Theorem (Functional description)

With a bounded psd kernel k , the next two claims are equivalent

1. k is characteristic
2. $\mathcal{H} + \mathbb{R}$ dense in all $L^2(R)$, for any probability measure R

Ex:

Gaussian and Laplace kernels

→ Difficult to check in general!

Theorem (Translation-invariant kernels)

- ▶ k : bounded psd kernel on \mathbb{R}^d with $k(x, y) = h(x - y)$
- ▶ $h(\cdot)$: bounded, continuous, positive and definite on \mathbb{R}^d

Then Bochner's theorem yields that h is the Fourier transform of a finite non-negative Borel measure Λ that is,

$$h(t) = \int_{\mathbb{R}^d} e^{-i\langle t, w \rangle} d\Lambda(w).$$

Moreover, the next two claims are equivalent

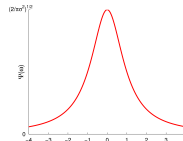
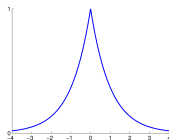
1. k is characteristic
2. $\text{Supp}(\Lambda) = \mathbb{R}^d$

Ex:

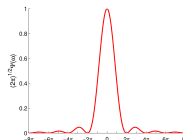
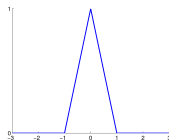
Gaussian and Laplace, Matern class, B_{2n+1} -splines, ... kernels

Characteristic kernels

- ▶ Gaussian: $h(t) = e^{-t^2/(2\sigma^2)}$, $\mathcal{F}(h)(u) = \sigma e^{-\sigma^2 u^2/2}$
- ▶ Laplace: $h(t) = e^{-|t|\sigma}$, $\mathcal{F}(h)(u) = \sqrt{2/\pi} \frac{\sigma}{\sigma^2 + u^2}$



- ▶ B_1 -Spline kernel: $h(t) = (1 - |t|)\mathbb{1}_{[-1,1]}(t)$,
 $\mathcal{F}(h)(u) = \frac{2\sqrt{2}}{\sqrt{\pi}} \frac{\sin(u/2)^2}{u^2}$

Distance between
Probability
distributionsMaximum Mean
Discrepancy
(MMD)

Mean element

Examples of mean
elements

Characteristic kernel

Two-sample test

Non characteristic kernels: Examples

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

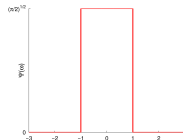
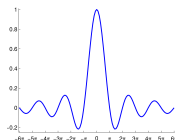
Examples of mean
elements

Characteristic kernel

Two-sample test

Non characteristic kernels

- Sinc kernel: $h(t) = \frac{\sin(\sigma t)}{t}$, $\mathcal{F}(h)(u) = \sqrt{\frac{\pi}{2}} \mathbb{1}_{[-\sigma, \sigma]}(u)$



- Periodic functions h on \mathbb{R}^d in full generality

Characteristic kernel

With a characteristic kernel,

$$P = Q \Leftrightarrow \mu_P = \mu_Q$$

MMD as a distance over probability measures

With a characteristic kernel,

$$P = Q \Leftrightarrow \text{MMD}_k(P, Q) = \|\mu_P - \mu_Q\|_{\mathcal{H}} = 0$$

Conclusion:

The two claims are equivalent:

- ▶ MMD_k : distance over probability distributions
- ▶ k : characteristic

Two-sample test

Testing for two fixed populations

Kernel Machines

Alain Celisse

Two-sample test

$X_1, \dots, X_n \sim P$ and $Y_1, \dots, Y_m \sim Q$, with $P = Q$?

Statistical test

- ▶ Null hypothesis:

$$H_0 : P = Q \text{ (no change)}$$

- ▶ Alternative:

$$H_1 : P \neq Q$$

- ▶ Oracle Test Statistic and Rejection region:

$$T = MMD_k(P, Q) = \|\mu_P - \mu_Q\|_{\mathcal{H}}$$

$$\mathcal{R} = \{MMD_k(P, Q) > 0\}$$

- ▶ Should be zero under H_0 and away from 0 otherwise
- ▶ To be estimated...

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Estimating the mean elements

$$\mu_P = \mathbb{E}[k(X, \cdot)] \approx \hat{\mu}_P = \frac{1}{n} \sum_{i=1}^n k(X_i, \cdot)$$

$$\mu_Q = \mathbb{E}[k(Y, \cdot)] \approx \hat{\mu}_Q = \frac{1}{m} \sum_{j=1}^m k(Y_j, \cdot)$$

First plug-in estimator of the MMD

$$\begin{aligned} \widehat{MMD}_{b,k}^2(P, Q) &= \|\hat{\mu}_P - \hat{\mu}_Q\|_{\mathcal{H}}^2 \\ &= \|\hat{\mu}_P\|_{\mathcal{H}}^2 + \|\hat{\mu}_Q\|_{\mathcal{H}}^2 - 2\langle \hat{\mu}_P, \hat{\mu}_Q \rangle_{\mathcal{H}} \\ &= \frac{\sum_{i,j=1}^n k(X_i, X_j)}{n^2} + \frac{\sum_{p,q=1}^m k(Y_p, Y_q)}{m^2} - 2 \frac{\sum_{i=1}^n \sum_{j=1}^m k(X_i, Y_j)}{mn} \end{aligned}$$

Distance between
Probability
distributionsMaximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distributionSingle change-point
detection

Independence testing

Bias of the plug-in estimator

Covariance operator

Definition

The covariance operator associated with X is the unique operator Σ^X from \mathcal{H} to \mathcal{H} such that

$$\langle \Sigma^X f, g \rangle_{\mathcal{H}} = \mathbb{E}[f(X) \cdot g(X)], \quad \forall f, g \in \mathcal{H}$$

Remark:

From the reproducing property:

$$\begin{aligned} \mathbb{E}[\langle f, k_X \rangle_{\mathcal{H}} \langle g, k_X \rangle_{\mathcal{H}}] &= \mathbb{E}[\langle (k_X \otimes k_X) f, g \rangle_{\mathcal{H}}] \\ &= \langle \mathbb{E}[k_X \otimes k_X] f, g \rangle_{\mathcal{H}} \end{aligned}$$

Theorem

$$\begin{aligned} &\mathbb{E} \left[\widehat{MMD}_{b,k}^2(P, Q) \right] \\ &= MMD_k^2(P, Q) + \frac{\text{Tr}(\Sigma^X) - \|\mu_P\|_{\mathcal{H}}^2}{n} + \frac{\text{Tr}(\Sigma^Y) - \|\mu_Q\|_{\mathcal{H}}^2}{n} \end{aligned}$$

Theorem

- \mathcal{K} : set of candidate kernels
- k : kernel such that $\sup_{k \in \mathcal{K}} \|k_x\|_{\mathcal{H}} \leq B$, for all $x \in \mathcal{X}$

With proba at least $1 - \delta$,

$$\begin{aligned} & \left| \widehat{MMD}_{b,k}(P, Q) - MMD_k(P, Q) \right| \\ & \leq C \left(\sqrt{\frac{\mathcal{Rad}_m(\mathcal{K})}{m}} + \sqrt{\frac{\mathcal{Rad}_n(\mathcal{K})}{n}} \right) \\ & \quad + C' B \left(1 + \sqrt{\log \left(\frac{4}{\delta} \right)} \right) \cdot \sqrt{\frac{m+n}{mn}} \end{aligned}$$

where the empirical Rademacher complexity is given by

$$\mathcal{Rad}_n(\mathcal{K}) = \mathbb{E} \left[\sup_{k \in \mathcal{K}} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j k(X_i, X_j) \right| \mid X_1, \dots, X_n \right]$$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Designing an estimator: 2nd try ...

$$\begin{aligned}MMD_k^2(P, Q) &= \|\mu_P\|_{\mathcal{H}}^2 + \|\mu_Q\|_{\mathcal{H}}^2 - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\&= \langle \mu_P, \mu_P \rangle_{\mathcal{H}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{H}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\&= E_{X, X'} [\langle k_X, k_{X'} \rangle_{\mathcal{H}}] + E_{Y, Y'} [\langle k_Y, k_{Y'} \rangle_{\mathcal{H}}] - 2 E_{X, Y} [\langle k_X, k_Y \rangle_{\mathcal{H}}]\end{aligned}$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Designing an estimator: 2nd try ...

$$\begin{aligned} \text{MMD}_k^2(P, Q) &= \|\mu_P\|_{\mathcal{H}}^2 + \|\mu_Q\|_{\mathcal{H}}^2 - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{H}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{H}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\ &= E_{X, X'} [\langle k_X, k_{X'} \rangle_{\mathcal{H}}] + E_{Y, Y'} [\langle k_Y, k_{Y'} \rangle_{\mathcal{H}}] - 2 E_{X, Y} [\langle k_X, k_Y \rangle_{\mathcal{H}}] \end{aligned}$$

Second plug-in estimator of the MMD

$$\begin{aligned} \widehat{\text{MMD}}_{u,k}^2(P, Q) &= \frac{1}{n(n-1)} \sum_{i \neq j=1}^n k(X_i, X_j) + \frac{1}{m(m-1)} \sum_{p \neq q=1}^m k(Y_p, Y_q) \\ &\quad - 2 \frac{\sum_{i=1}^n \sum_{j=1}^m k(X_i, Y_j)}{mn} \end{aligned}$$

Designing an estimator: 2nd try ...

$$\begin{aligned}MMD_k^2(P, Q) &= \|\mu_P\|_{\mathcal{H}}^2 + \|\mu_Q\|_{\mathcal{H}}^2 - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\&= \langle \mu_P, \mu_P \rangle_{\mathcal{H}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{H}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\&= E_{X, X'} [\langle k_X, k_{X'} \rangle_{\mathcal{H}}] + E_{Y, Y'} [\langle k_Y, k_{Y'} \rangle_{\mathcal{H}}] - 2 E_{X, Y} [\langle k_X, k_Y \rangle_{\mathcal{H}}]\end{aligned}$$

Second plug-in estimator of the MMD

$$\begin{aligned}\widehat{MMD}_{u,k}^2(P, Q) \\&= \frac{1}{n(n-1)} \sum_{i \neq j=1}^n k(X_i, X_j) + \frac{1}{m(m-1)} \sum_{p \neq q=1}^m k(Y_p, Y_q) \\&\quad - 2 \frac{\sum_{i=1}^n \sum_{j=1}^m k(X_i, Y_j)}{mn}\end{aligned}$$

Theorem

$$\mathbb{E} \left[\widehat{MMD}_{u,k}^2(P, Q) \right] = MMD_k^2(P, Q)$$

Distance between
Probability
distributionsMaximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distributionSingle change-point
detection

Independence testing

Designing an estimator: 2nd try ...

$$\begin{aligned} \text{MMD}_k^2(P, Q) &= \|\mu_P\|_{\mathcal{H}}^2 + \|\mu_Q\|_{\mathcal{H}}^2 - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{H}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{H}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{H}} \\ &= E_{X, X'} [\langle k_X, k_{X'} \rangle_{\mathcal{H}}] + E_{Y, Y'} [\langle k_Y, k_{Y'} \rangle_{\mathcal{H}}] - 2 E_{X, Y} [\langle k_X, k_Y \rangle_{\mathcal{H}}] \end{aligned}$$

Second plug-in estimator of the MMD

$$\begin{aligned} \widehat{\text{MMD}}_{u,k}^2(P, Q) &= \frac{1}{n(n-1)} \sum_{i \neq j=1}^n k(X_i, X_j) + \frac{1}{m(m-1)} \sum_{p \neq q=1}^m k(Y_p, Y_q) \\ &\quad - 2 \frac{\sum_{i=1}^n \sum_{j=1}^m k(X_i, Y_j)}{mn} \end{aligned}$$

Theorem

$$\mathbb{E} \left[\widehat{\text{MMD}}_{u,k}^2(P, Q) \right] = \text{MMD}_k^2(P, Q)$$

Remark: Is there another choice for an unbiased estimator?

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Testing for two fixed populations

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Two-sample test

$X_1, \dots, X_n \sim P$ and $Y_1, \dots, Y_m \sim Q$, with $P = Q$?

Statistical test

► Hypothesis:

$$H_0 : P = Q \text{ (no change)} \quad \text{vs} \quad H_1 : P \neq Q$$

► Rejection region:

$$\mathcal{R} = \left\{ \widehat{MMD}_k^2(P, Q) > \eta_\alpha \right\}$$

where $\eta_\alpha > 0$ depends on the Type-I error

Problem

Distribution of $\widehat{MMD}_k^2(P, Q)$ difficult to estimate

Distribution of the unbiased estimator

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

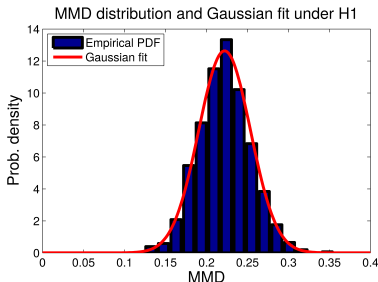
Single change-point
detection

Independence testing

$$\begin{aligned} & \widehat{MMD}_{u,k}^2(P, Q) \\ &= \frac{1}{n(n-1)} \sum_{i \neq j=1}^n k(X_i, X_j) + \frac{1}{m(m-1)} \sum_{p \neq q=1}^m k(Y_p, Y_q) \\ & \quad - 2 \frac{\sum_{i=1}^n \sum_{j=1}^m k(X_i, Y_j)}{mn} \end{aligned}$$

If $P \neq Q$ ($m = n$)

$\widehat{MMD}_{u,k}^2(P, Q)$: Gaussian asymptotic distrib. (U -stat.)



Distribution of the unbiased estimator

Kernel Machines

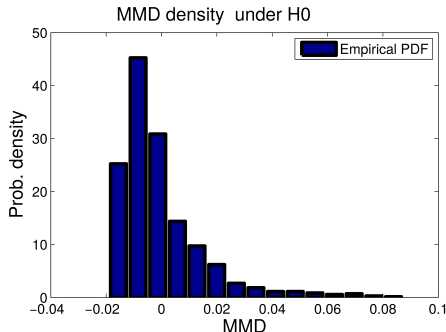
Alain Celisse

If $P = Q$ ($m = n$)

$$m\widehat{MMD}_k^2(P, Q) \sim 2 \sum_{\ell=1}^{+\infty} \lambda_{\ell} [\chi_{\ell}^2 - 1]$$

where, for all $\ell \geq 1$,

$$\int_{\mathcal{X}} k_x \psi_{\ell}(x) dP(x) = \lambda_{\ell} \psi_{\ell} \in \mathcal{H}$$



Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Sequential testing for one change-point

Kernel Machines

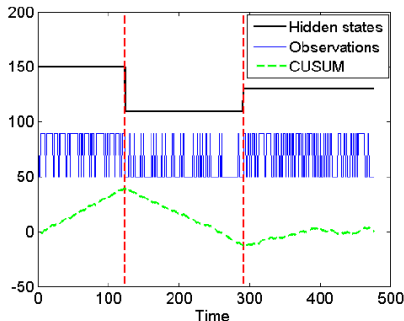
Alain Celisse

Single changepoint detection

- ▶ Time series: $(t_1, X_1), \dots, (t_n, X_n) \in [0, +\infty] \times \mathbb{R}$
- ▶ Assumption:

$$P_{X_1} = \dots = P_{X_r} \neq P_{X_{r+1}} = \dots = P_{X_n}$$

- ▶ Usually replaced by the mean (distributional features)



Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

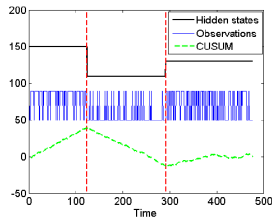
MMD estimator
distribution

Single change-point
detection

Independence testing

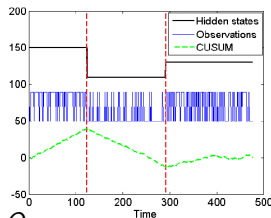
From two-sample test to Single changepoint detection

- ▶ r : changepoint location
- ▶ $P_{X_r} \neq P_{X_{r+1}}$: “abrupt” change



From two-sample test to Single changepoint detection

- ▶ r : changepoint location
- ▶ $P_{X_r} \neq P_{X_{r+1}}$: “abrupt” change
- ▶ Two populations:
 - ▶ From 1 to r : all X_i s from P
 - ▶ From $r + 1$ to n : all X_i s from Q
- ▶ Each time $1 \leq t \leq n$: candidate changepoint



Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

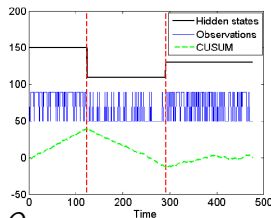
Independence testing

From two-sample test to Single changepoint detection

Kernel Machines

Alain Celisse

- ▶ r : changepoint location
- ▶ $P_{X_r} \neq P_{X_{r+1}}$: “abrupt” change
- ▶ Two populations:
 - ▶ From 1 to r : all X_i s from P
 - ▶ From $r + 1$ to n : all X_i s from Q
- ▶ Each time $1 \leq t \leq n$: candidate changepoint



MMD-based statistic

$$\begin{aligned}\widehat{MMD}_{u,k}^2(t) &= \|\widehat{\mu}_P(1:t) - \widehat{\mu}_Q(t+1:n)\|_{\mathcal{H}}^2 \\ &= \frac{\sum_{i \neq j=1}^t k(X_i, X_j)}{t(t-1)} + \frac{\sum_{p \neq q=t+1}^n k(Y_p, Y_q)}{(n-t)(n-t-1)} \\ &\quad - 2 \frac{\sum_{i=1}^t \sum_{j=t+1}^n k(X_i, Y_j)}{t(n-t)}\end{aligned}$$

Distance between Probability distributions

Maximum Mean Discrepancy (MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator distribution

Single change-point detection

Independence testing

Single changepoint detection algorithm

Algorithm

1. Compute

$$\mathcal{M}_n = \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Single changepoint detection algorithm

Algorithm

1. Compute

$$\mathcal{M}_n = \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

2. Rejection region:

$$\mathcal{R}_\alpha = \{ \mathcal{M}_n \geq q_\alpha \}$$

q_α : α -quantile of $P_{\mathcal{M}_n}$ under the null

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Single changepoint detection algorithm

Kernel Machines

Alain Celisse

Algorithm

1. Compute

$$\mathcal{M}_n = \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

2. Rejection region:

$$\mathcal{R}_\alpha = \{ \mathcal{M}_n \geq q_\alpha \}$$

q_α : α -quantile of $P_{\mathcal{M}_n}$ under the null

3. If the null hypothesis is rejected, then

$$\hat{t}_\alpha = \text{Arg} \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Single changepoint detection algorithm

Kernel Machines

Alain Celisse

Algorithm

1. Compute

$$\mathcal{M}_n = \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

2. Rejection region:

$$\mathcal{R}_\alpha = \{ \mathcal{M}_n \geq q_\alpha \}$$

q_α : α -quantile of $P_{\mathcal{M}_n}$ under the null

3. If the null hypothesis is rejected, then

$$\hat{r}_\alpha = \text{Arg} \min_{1 \leq t \leq n} \left\{ \widehat{MMD}_{u,k}^2(t) \right\}$$

Remarks:

- ▶ q_α computed by bootstrap or approximation (asymptotic Gaussian)
- ▶ Sequential approaches also possible (multiple changepoints)

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Independent random variables?

Kernel Machines

Alain Celisse

Two samples

- ▶ $(X_i, Y_i)_{i=1}^n$: n couples
- ▶ $X_1, \dots, X_n \sim P_X$: realizations of X
- ▶ $Y_1, \dots, Y_n \sim P_Y$: realizations of Y

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Independent random variables?

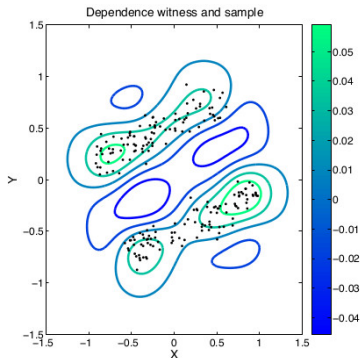
Kernel Machines

Alain Celisse

Two samples

- ▶ $(X_i, Y_i)_{i=1}^n$: n couples
- ▶ $X_1, \dots, X_n \sim P_X$: realizations of X
- ▶ $Y_1, \dots, Y_n \sim P_Y$: realizations of Y

Are X and Y independent random variables?



Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

Statistical test

- ▶ Null hypothesis:

$$H_0 : P_{X,Y} = P_X \otimes P_Y$$

- ▶ Alternative:

$$H_1 : P_{X,Y} \neq P_X \otimes P_Y$$

Statistical test

- ▶ Null hypothesis:

$$H_0 : P_{X,Y} = P_X \otimes P_Y$$

- ▶ Alternative:

$$H_1 : P_{X,Y} \neq P_X \otimes P_Y$$

Idea

- ▶ $P = P_{X,Y}$
- ▶ $Q = P_X \otimes P_Y$

Statistical test

- ▶ Null hypothesis:

$$H_0 : P_{X,Y} = P_X \otimes P_Y$$

- ▶ Alternative:

$$H_1 : P_{X,Y} \neq P_X \otimes P_Y$$

Idea

- ▶ $P = P_{X,Y}$
- ▶ $Q = P_X \otimes P_Y$

$$MMD_k(P, Q) = \|\mu_P - \mu_Q\|_{\mathcal{H}_X, \mathcal{H}_Y}^2$$

HSIC criterion for independence testing

Product kernel on $\mathcal{X} \times \mathcal{Y}$

$$k((X_i, Y_i), (X_j, Y_j)) = k_X(X_i, X_j) \times k_Y(Y_i, Y_j)$$

- ▶ \mathcal{H}_X : RKHS of k_X
- ▶ \mathcal{H}_Y : RKHS of k_Y

Kernel Machines

Alain Celisse

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing

HSIC criterion for independence testing

Product kernel on $\mathcal{X} \times \mathcal{Y}$

$$k((X_i, Y_i), (X_j, Y_j)) = k_X(X_i, X_j) \times k_Y(Y_i, Y_j)$$

- ▶ \mathcal{H}_X : RKHS of k_X
- ▶ \mathcal{H}_Y : RKHS of k_Y
- ▶ Then: $k(\cdot, \cdot)$: reproducing kernel on $\mathcal{H}_X \times \mathcal{H}_Y$

Remark:

Other choices of kernels are possible

HSIC

$$\begin{aligned} \text{HSIC}(P_{X,Y}, P_X \otimes P_Y) &= \|\mu_{P_{X,Y}} - \mu_{P_X \otimes P_Y}\|_{\mathcal{H}_X, \mathcal{H}_Y}^2 \\ &= \text{MMD}_k^2(P_{X,Y}, P_X \otimes P_Y) \end{aligned}$$

HSIC criterion for independence testing

Kernel Machines

Alain Celisse

Product kernel on $\mathcal{X} \times \mathcal{Y}$

$$k((X_i, Y_i), (X_j, Y_j)) = k_X(X_i, X_j) \times k_Y(Y_i, Y_j)$$

- ▶ \mathcal{H}_X : RKHS of k_X
- ▶ \mathcal{H}_Y : RKHS of k_Y
- ▶ Then: $k(\cdot, \cdot)$: reproducing kernel on $\mathcal{H}_X \times \mathcal{H}_Y$

Remark:

Other choices of kernels are possible

HSIC

$$\begin{aligned} \text{HSIC}(P_{X,Y}, P_X \otimes P_Y) &= \|\mu_{P_{X,Y}} - \mu_{P_X \otimes P_Y}\|_{\mathcal{H}_X, \mathcal{H}_Y}^2 \\ &= \text{MMD}_k^2(P_{X,Y}, P_X \otimes P_Y) \end{aligned}$$

<w**Exercise:** Compute $\|\mu_{P_{X,Y}}\|_{\mathcal{H}_X \times \mathcal{H}_Y}^2$ and $\|\mu_{P_X \otimes P_Y}\|_{\mathcal{H}_X \times \mathcal{H}_Y}^2$

Distance between
Probability
distributions

Maximum Mean
Discrepancy
(MMD)

Mean element

Two-sample test

Setting

Estimating the MMD

MMD estimator
distribution

Single change-point
detection

Independence testing