Modern Kernel Methods in Machine Learning: Part II

Danica J. Sutherland (she/her) Computer Science, University of British Columbia ETICS "summer" school, Oct 2022

Yesterday, we saw:

- RKHS \mathcal{H} is a function space, $f:\mathcal{X}
 ightarrow\mathbb{R}$
- Reproducing property: $\langle f, k(x, \cdot)
 angle_{\mathcal{H}} = f(x)$
- Representer theorem: $rgmin L(f(x_1),\ldots,f(x_n))+R(\|f\|_{\mathcal{H}})\in \mathrm{span}\{k(x_i,\cdot)\}_{i=1}^n$
- Can use to do kernel ridge regression, SVMs, etc

- Kernel mean embeddings of distributions
- Gaussian processes and probabilistic numerics
- Kernel approximations, for better computation
- Neural tangent kernels

- Kernel mean embeddings of distributions
- Gaussian processes and probabilistic numerics
- Kernel approximations, for better computation
- Neural tangent kernels

- Kernel mean embeddings of distributions
- Gaussian processes and probabilistic numerics
- Kernel approximations, for better computation
- Neural tangent kernels

- Kernel mean embeddings of distributions
- Gaussian processes and probabilistic numerics
- Kernel approximations, for better computation
- Neural tangent kernels

• Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim \mathbb P} f(X) = \langle f, \mu_{\mathbb P}
 angle_{\mathcal H}$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

 $\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} | k(X, \cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

 $\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} | k(X, \cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

• Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} | k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$
- $\bullet \ \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} = \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \ k(X, Y)$

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$

- $\bullet \ \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} = \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \ k(X, Y)$
- Okay. Why?

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$

- $\bullet \ \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} = \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \ k(X, Y)$
- Okay. Why?
 - One reason: ML on distributions [Szabó+ JMLR-16]

- Represent point $x \in \mathcal{X}$ as $k(x, \cdot)$: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$: $\mathbb E_{X\sim\mathbb P} f(X) = \langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$

- $\bullet \ \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} = \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \ k(X, Y)$
- Okay. Why?
 - One reason: ML on distributions [Szabó+ JMLR-16]
 - More common reason: comparing distributions

 $\mathrm{MMD}(\mathbb{P},\mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}}$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|oldsymbol{\mu}_{\mathbb{P}} - oldsymbol{\mu}_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, oldsymbol{\mu}_{\mathbb{P}} - oldsymbol{\mu}_{\mathbb{Q}}
angle_{\mathcal{H}} \end{aligned}$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} f(Y) \end{aligned}$$

• Last line is Integral Probability Metric (IPM) form

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \: k(t, X) - \mathbb{E}_{\mathbb{Q}} \: k(t, Y)$$



$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^*(t) \propto \langle \mu_\mathbb{P} - \mu_\mathbb{Q}, k(t, \cdot)
angle_\mathcal{H} = \mathbb{E}_\mathbb{P} \: k(t, X) - \mathbb{E}_\mathbb{Q} \: k(t, Y)$$



$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^{*}(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^{*}(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^{*}(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \| \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on \mathbb{P} , low on \mathbb{Q}

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

• $MMD(\mathbb{P},\mathbb{P}) = 0$, symmetry, triangle inequality

- $\mathrm{MMD}(\mathbb{P},\mathbb{P})=0$, symmetry, triangle inequality
- If k is *characteristic*, then $MMD(\mathbb{P}, \mathbb{Q}) = 0$ iff $\mathbb{P} = \mathbb{Q}$ • i.e. $\mathbb{P} \mapsto \mu_{\mathbb{P}}$ is injective

- $\mathrm{MMD}(\mathbb{P},\mathbb{P})=0$, symmetry, triangle inequality
- If k is *characteristic*, then $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
 - i.e. $\mathbb{P}\mapsto \mu_{\mathbb{P}}$ is injective
 - Makes MMD a metric on probability distributions

- $\mathrm{MMD}(\mathbb{P},\mathbb{P})=0$, symmetry, triangle inequality
- If k is *characteristic*, then $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
 - i.e. $\mathbb{P}\mapsto \mu_{\mathbb{P}}$ is injective
 - Makes MMD a metric on probability distributions
 - Universal => characteristic

- $\mathrm{MMD}(\mathbb{P},\mathbb{P})=0$, symmetry, triangle inequality
- If k is *characteristic*, then $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
 - i.e. $\mathbb{P}\mapsto \mu_{\mathbb{P}}$ is injective
 - Makes MMD a metric on probability distributions
 - Universal => characteristic
- Linear kernel: $\mathrm{MMD}(\mathbb{P},\mathbb{Q}) = \|\mu_{\mathbb{P}} \mu_{\mathbb{Q}}\|_{\mathcal{H}}$ is just Euclidean distance between means

• Want a "super-sample" from $\mathbb{P}: rac{1}{n}\sum_i f(X_i) pprox \mathbb{E}\,f(X)$

- Want a "super-sample" from \mathbb{P} : $\frac{1}{n} \sum_i f(X_i) \approx \mathbb{E} f(X)$
- If $f \in \mathcal{H}$, error $\leq \|f\|_{\mathcal{H}} \operatorname{MMD}(\mathbb{P}, rac{1}{n} \sum_{i=1}^T \delta_{X_i})$

• Want a "super-sample" from \mathbb{P} : $rac{1}{n}\sum_i f(X_i) pprox \mathbb{E}\,f(X)$

- If $f \in \mathcal{H}$, error $\leq \|f\|_{\mathcal{H}} \operatorname{MMD}(\mathbb{P}, rac{1}{n} \sum_{i=1}^T \delta_{X_i})$
- Greedily minimize the MMD:

 $X_{T+1} \in rgmin_{X \in \mathcal{X}} \mathbb{E}_{X' \sim \mathbb{P}} \, k(X,X') - rac{1}{T+1} \sum_{i=1}^T k(X,X_i)$

• Want a "super-sample" from \mathbb{P} : $rac{1}{n}\sum_i f(X_i) pprox \mathbb{E}\,f(X)$

- If $f \in \mathcal{H}$, error $\leq \|f\|_{\mathcal{H}} \operatorname{MMD}(\mathbb{P}, rac{1}{n} \sum_{i=1}^T \delta_{X_i})$
- Greedily minimize the MMD:

$$X_{T+1} \in rgmin_{X \in \mathcal{X}} \mathbb{E}_{X' \sim \mathbb{P}} k(X, X') - rac{1}{T+1} \sum_{i=1}^T k(X, X_i)$$

• Get $\mathcal{O}(1/T)$ approximation instead of $\mathcal{O}(1/\sqrt{T})$ with random samples


$\mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) = \mathbb{E}_{\substack{X,X'\sim\mathbb{P}\Y,Y'\sim\mathbb{Q}}}\left[k(X,X')-2k(X,Y)+k(Y,Y')
ight]$

 $\mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) = \mathbb{E}_{\substack{X,X'\sim\mathbb{P}\Y,Y'\sim\mathbb{Q}}}\left[k(X,X')-2k(X,Y)+k(Y,Y')
ight]$

 $\widehat{\mathrm{MMD}}_k^2(X, Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \mathrm{mean}(K_{XY})$

 $\mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) = \mathbb{E}_{\substack{X,X'\sim\mathbb{P}\Y,Y'\sim\mathbb{Q}}} \left[k(X,X') - 2k(X,Y) + k(Y,Y')
ight]$

 $\widehat{\mathrm{MMD}}_k^2(X,Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \mathrm{mean}(K_{XY})$

K_{XX}



 $\mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) = \mathbb{E}_{\substack{X,X'\sim\mathbb{P}\Y,Y'\sim\mathbb{Q}}} \left[k(X,X') - 2k(X,Y) + k(Y,Y')
ight]$

 $\widehat{\mathrm{MMD}}_k^2(X,Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \mathrm{mean}(K_{XY})$

 K_{XX}

 K_{YY}

1.0	0.2	0.6		1.0	0.8	0.7
0.2	1.0	0.5	· <u>(</u>),	0.8	1.0	0.6
0.6	0.5	1.0		0.7	0.6	1.0

 $\mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) = \mathbb{E}_{\substack{X,X'\sim\mathbb{P}\Y,Y'\sim\mathbb{Q}}} \left[k(X,X') - 2k(X,Y) + k(Y,Y')
ight]$

 $\widehat{\mathrm{MMD}}_k^2(X,Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \, \mathrm{mean}(K_{XY})$



1.0	0.2	0.6		1.0	0.8	0.7	 0.1	0.2
0.2	1.0	0.5		0.8	1.0	0.6	 0.3	0.3
0.6	0.5	1.0	(Cinin _ ami);	0.7	0.6	1.0	 0.1	0.4

- MMD has easy $\mathcal{O}(n^2)$ estimator
 - block or incomplete estimators are $\mathcal{O}(n^{lpha})$ for $lpha \in [1,2]$, but noisier

- MMD has easy $\mathcal{O}(n^2)$ estimator
 - block or incomplete estimators are $\mathcal{O}(n^{lpha})$ for $lpha \in [1,2]$, but noisier
- For bounded kernel, $\mathcal{O}_p(1/\sqrt{n})$ estimation error

- MMD has easy $\mathcal{O}(n^2)$ estimator
 - block or incomplete estimators are $\mathcal{O}(n^{lpha})$ for $lpha \in [1,2]$, but noisier
- For bounded kernel, $\mathcal{O}_p(1/\sqrt{n})$ estimation error
 - Independent of data dimension!

- MMD has easy $\mathcal{O}(n^2)$ estimator
 - block or incomplete estimators are $\mathcal{O}(n^{lpha})$ for $lpha \in [1,2]$, but noisier
- For bounded kernel, $\mathcal{O}_p(1/\sqrt{n})$ estimation error
 - Independent of data dimension!
 - But, no free lunch...the *value* of the MMD generally shrinks with growing dimension, so constant $\mathcal{O}_p(1/\sqrt{n})$ error gets worse relatively



• Given samples from two unknown distributions

 $X \sim \mathbb{P}$ $Y \sim \mathbb{Q}$

• Question: is $\mathbb{P} = \mathbb{Q}$?

• Given samples from two unknown distributions

 $X\sim \mathbb{P} \qquad Y\sim \mathbb{O}$

• Do smokers/non-smokers get different cancers?

• Given samples from two unknown distributions

 $X\sim \mathbb{P} \qquad Y\sim \mathbb{O}$

- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?

- Given samples from two unknown distributions
 - $X \sim \mathbb{P}$ $Y \sim \mathbb{O}$
- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?
- When does my laser agree with the one on Mars?

• Given samples from two unknown distributions

- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?
- When does my laser agree with the one on Mars?
- Are storms in the 2000s different from storms in the 1800s?

• Given samples from two unknown distributions

- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?
- When does my laser agree with the one on Mars?
- Are storms in the 2000s different from storms in the 1800s?
- Does presence of this protein affect DNA binding? [MMDiff2]

• Given samples from two unknown distributions

- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?
- When does my laser agree with the one on Mars?
- Are storms in the 2000s different from storms in the 1800s?
- Does presence of this protein affect DNA binding? [MMDiff2]
- Do these dob and birthday columns mean the same thing?

• Given samples from two unknown distributions

- Do smokers/non-smokers get different cancers?
- Do Brits have the same friend network types as Americans?
- When does my laser agree with the one on Mars?
- Are storms in the 2000s different from storms in the 1800s?
- Does presence of this protein affect DNA binding? [MMDiff2]
- Do these dob and birthday columns mean the same thing?
- Does my generative model \mathbb{Q}_{θ} match \mathbb{P}_{data} ?

• Given samples from two unknown distributions

 $X \sim \mathbb{P}$ $Y \sim \mathbb{Q}$

• Question: is $\mathbb{P} = \mathbb{Q}$?

• Given samples from two unknown distributions

 $X \sim \mathbb{P}$ $Y \sim \mathbb{Q}$

- Question: is $\mathbb{P} = \mathbb{Q}$?
- Hypothesis testing approach:

$$H_0:\mathbb{P}=\mathbb{Q} \qquad H_1:\mathbb{P}
eq \mathbb{Q}$$

• Given samples from two unknown distributions

 $X \sim \mathbb{P}$ $Y \sim \mathbb{Q}$

- Question: is $\mathbb{P} = \mathbb{Q}$?
- Hypothesis testing approach:

 $H_0:\mathbb{P}=\mathbb{Q} \qquad H_1:\mathbb{P}
eq \mathbb{Q}$

• Reject $\widehat{H_0}$ if $\widehat{\mathrm{MMD}}(X,Y) > c_lpha$











• $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something

Infinite mixture of χ^2 s, params depend on \mathbb{P} and k

- $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something
 - Infinite mixture of χ^2 s, params depend on \mathbb{P} and k
 - Can estimate threshold with *permutation testing*

- $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something
 - Infinite mixture of χ^2 s, params depend on \mathbb{P} and k
 - Can estimate threshold with *permutation testing*

•
$$H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 - \mathrm{MMD}^2) \overset{d}{
ightarrow}$$
 asymptotically normal

- $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something
 - Infinite mixture of χ^2 s, params depend on \mathbb{P} and k
 - Can estimate threshold with permutation testing
- $H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 \mathrm{MMD}^2) \overset{d}{
 ightarrow}$ asymptotically normal
- Any characteristic kernel gives consistent test

- $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something
 - Infinite mixture of χ^2 s, params depend on \mathbb{P} and k
 - Can estimate threshold with permutation testing

•
$$H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 - \mathrm{MMD}^2) \stackrel{d}{
ightarrow}$$
 asymptotically normal

• Any characteristic kernel gives consistent test...eventually

- $H_0: n \widehat{\mathrm{MMD}}^2$ converges in distribution to...something
 - Infinite mixture of χ^2 s, params depend on \mathbb{P} and k
 - Can estimate threshold with permutation testing

•
$$H_1: \sqrt{n} (\widehat{\mathrm{MMD}}^2 - \mathrm{MMD}^2) \overset{d}{
ightarrow}$$
 asymptotically normal

- Any characteristic kernel gives consistent test...eventually
- Need enormous $oldsymbol{n}$ if kernel is bad for problem

Classifier two-sample tests



- $\hat{T}(X, Y)$ is the accuracy of f on the test set
- Under H_0 , classification impossible: $\hat{T} \sim \mathrm{Binomial}(n, rac{1}{2})$

Classifier two-sample tests



- $\hat{T}(X, Y)$ is the accuracy of f on the test set
- Under H_0 , classification impossible: $\hat{T} \sim \mathrm{Binomial}(n, rac{1}{2})$
- With $k(x,y)=rac{1}{4}f(x)f(y)$ where $f(x)\in\{-1,1\}$, get $\widehat{\mathrm{MMD}}(X,Y)=\left|\hat{T}(X,Y)-rac{1}{2}
 ight|$

Deep learning and deep kernels

• $k(x,y) = rac{1}{4}f(x)f(y)$ is one form of *deep kernel*

Deep learning and deep kernels

- $k(x,y)=rac{1}{4}f(x)f(y)$ is one form of deep kernel
- Deep models are usually of the form $f(x) = w^{\mathsf{T}} \phi_\psi(x)$
 - With a learned $\phi_\psi(x):\mathcal{X} o\mathbb{R}^D$
Deep learning and deep kernels

- $k(x,y) = rac{1}{4}f(x)f(y)$ is one form of deep kernel
- Deep models are usually of the form $f(x) = w^{\mathsf{T}} \phi_\psi(x)$
 - With a learned $\phi_\psi(x):\mathcal{X} o\mathbb{R}^D$
- If we fix ψ , have $f\in \mathcal{H}_\psi$ with $k_\psi(x,y)=\phi_\psi(x)^{\sf T}\phi_\psi(y)$

Deep learning and deep kernels

• $k(x,y) = rac{1}{4}f(x)f(y)$ is one form of deep kernel

- Deep models are usually of the form $f(x) = w^{\mathsf{T}} \phi_\psi(x)$
 - With a learned $\phi_\psi(x):\mathcal{X} o\mathbb{R}^D$
- If we fix ψ , have $f\in \mathcal{H}_\psi$ with $k_\psi(x,y)=\phi_\psi(x)^{\sf T}\phi_\psi(y)$
 - Same idea as NNGP approximation

Deep learning and deep kernels

- $k(x,y) = rac{1}{4}f(x)f(y)$ is one form of deep kernel
- Deep models are usually of the form $f(x) = w^{\mathsf{T}} \phi_\psi(x)$
 - With a learned $\phi_\psi(x):\mathcal{X} o\mathbb{R}^D$
- If we fix ψ , have $f\in \mathcal{H}_\psi$ with $k_\psi(x,y)=\phi_\psi(x)^{\sf T}\phi_\psi(y)$
 - Same idea as NNGP approximation
- Generalize to a **deep kernel**:

$$k_\psi(x,y) = \kappa\left(\phi_\psi(x),\phi_\psi(y)
ight)$$

• Take
$$k_\psi(x,y) = rac{1}{4} f_\psi(x) f_\psi(y)$$

• Final function in \mathcal{H}_ψ will be $af_\psi(x)$

• Take
$$k_\psi(x,y) = rac{1}{4} f_\psi(x) f_\psi(y) + 1$$

• Final function in \mathcal{H}_ψ will be $af_\psi(x)+b$

• Take
$$k_\psi(x,y) = rac{1}{4} f_\psi(x) f_\psi(y) + 1$$

- Final function in \mathcal{H}_ψ will be $af_\psi(x)+b$
- With logistic loss: this is Platt scaling

• Take
$$k_\psi(x,y) = rac{1}{4} f_\psi(x) f_\psi(y) + 1$$

- Final function in \mathcal{H}_ψ will be $af_\psi(x)+b$
- With logistic loss: this is Platt scaling

On Calibration of Modern Neural Networks

Chuan Guo^{*1} **Geoff Pleiss**^{*1} **Yu Sun**^{*1} **Kilian Q. Weinberger**¹

• This definitely does *not* say that deep learning is (even approximately) a kernel method

- This definitely does *not* say that deep learning is (even approximately) a kernel method
- ...despite what some people might want you to think

Computer Science > Machine Learning

[Submitted on 30 Nov 2020]

Every Model Learned by Gradient Descent Is Approximately a Kernel Machine

Pedro Domingos

- This definitely does *not* say that deep learning is (even approximately) a kernel method
- ...despite what some people might want you to think

Computer Science > Machine Learning [Submitted on 30 Nov 2020] Every Model Learned by Gradient Descent Is Approximately a Kernel Machine Pedro Domingos

• We know theoretically deep learning can learn some things faster than any kernel method [see Malach+ ICML-21 + refs]

- This definitely does *not* say that deep learning is (even approximately) a kernel method
- ...despite what some people might want you to think

Computer Science > Machine Learning [Submitted on 30 Nov 2020] Every Model Learned by Gradient Descent Is Approximately a Kernel Machine Pedro Domingos

- We know theoretically deep learning can learn some things faster than any kernel method [see Malach+ ICML-21 + refs]
- But deep kernel learning ≠ traditional kernel models
 - exactly like how usual deep learning ≠ linear models

• Asymptotics of $\widehat{\mathrm{MMD}}^2$ give us immediately that

$$\Pr_{H_1}\left(n\widehat{ ext{MMD}}^2 > c_lpha
ight) pprox \Phi\left(rac{\sqrt{n}\, ext{MMD}^2}{\sigma_{H_1}} - rac{c_lpha}{\sqrt{n}\sigma_{H_1}}
ight)$$

- Asymptotics of \widehat{MMD}^2 give us immediately that

$$\Pr_{H_1}\left(n\widehat{ ext{MMD}}^2 > c_lpha
ight) pprox \Phi\left(rac{\sqrt{n}\, ext{MMD}^2}{\sigma_{H_1}} - rac{c_lpha}{\sqrt{n}\sigma_{H_1}}
ight)$$

 MMD , σ_{H_1} , c_lpha are constants: first term usually dominates

- Pick k to maximize an estimate of $\mathrm{MMD}^2 \, / \sigma_{H_1}$

• Asymptotics of $\widehat{\mathrm{MMD}}^2$ give us immediately that

$$\Pr_{H_1}\left(\widehat{ ext{MMD}}^2 > c_lpha
ight) pprox \Phi\left(rac{\sqrt{n}\, ext{MMD}^2}{\sigma_{H_1}} - rac{c_lpha}{\sqrt{n}\sigma_{H_1}}
ight)$$

- Pick k to maximize an estimate of $\mathrm{MMD}^2 \, / \sigma_{H_1}$
- Use $\widetilde{\mathrm{MMD}}$ from before, get $\hat{\sigma}_{H_1}$ from U-statistic theory

• Asymptotics of $\widehat{\mathrm{MMD}}^2$ give us immediately that

$$\Pr_{H_1}\left(n\widehat{ ext{MMD}}^2 > c_lpha
ight) pprox \Phi\left(rac{\sqrt{n}\, ext{MMD}^2}{\sigma_{H_1}} - rac{c_lpha}{\sqrt{n}\sigma_{H_1}}
ight)$$

- Pick k to maximize an estimate of $\mathrm{MMD}^2 \, / \sigma_{H_1}$
- Use $\widetilde{\mathrm{MMD}}$ from before, get $\hat{\sigma}_{H_1}$ from U-statistic theory
- Can show uniform $\mathcal{O}_P(n^{-\frac{1}{3}})$ convergence of estimator

• Asymptotics of $\widehat{\mathrm{MMD}}^2$ give us immediately that

$$\Pr_{H_1}\left(n\widehat{ ext{MMD}}^2 > c_lpha
ight) pprox \Phi\left(rac{\sqrt{n}\, ext{MMD}^2}{\sigma_{H_1}} - rac{c_lpha}{\sqrt{n}\sigma_{H_1}}
ight)$$

- Pick k to maximize an estimate of $\mathrm{MMD}^2 \, / \sigma_{H_1}$
- Use $\widehat{\mathrm{MMD}}$ from before, get $\hat{\sigma}_{H_1}$ from U-statistic theory
- Can show uniform $\mathcal{O}_P(n^{-\frac{1}{3}})$ convergence of estimator
- Get better tests (even after data splitting)

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}
 - Minimize estimate of $\operatorname{MMD}\psi(\mathbb{P}^m,\mathbb{Q}^n_{\theta})$ on a minibatch

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}
 - Minimize estimate of $\operatorname{MMD}\psi(\mathbb{P}^m,\mathbb{Q}^n_{\theta})$ on a minibatch
- MMD GAN: $\min_{\theta} [\max_{\psi} MMD_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta})]$

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}
 - Minimize estimate of $\operatorname{MMD}\psi(\mathbb{P}^m,\mathbb{Q}^n_{ heta})$ on a minibatch
- MMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{MMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
- SMMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{SMMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
 - Scaled MMD uses kernel properties to ensure smooth loss for θ by making witness function smooth [Arbel+ NeurIPS-18]

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}
 - Minimize estimate of $\operatorname{MMD}\psi(\mathbb{P}^m,\mathbb{Q}^n_{ heta})$ on a minibatch
- MMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{MMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
- SMMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{SMMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
 - Scaled MMD uses kernel properties to ensure smooth loss for θ by making witness function smooth [Arbel+ NeurIPS-18]
 - Uses $\langle f, \partial_{x_1} k(x, \cdot)
 angle_{\mathcal{H}} = \partial_{x_1} f(x)$

- An implicit generative model:
 - A generator net outputs samples from \mathbb{Q}_{θ}
 - Minimize estimate of $\operatorname{MMD}\psi(\mathbb{P}^m,\mathbb{Q}^n_{ heta})$ on a minibatch
- MMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{MMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
- SMMD GAN: $\min_{\theta} \left[\max_{\psi} \operatorname{SMMD}_{\psi}(\mathbb{P}, \mathbb{Q}_{\theta}) \right]$
 - Scaled MMD uses kernel properties to ensure smooth loss for θ by making witness function smooth [Arbel+ NeurIPS-18]
 - Uses $\langle f, \partial_{x_1} k(x, \cdot)
 angle_{\mathcal{H}} = \partial_{x_1} f(x)$
 - Standard WGAN-GP better thought of in kernel framework

Application: distribution regression/classification/...

• We can define a kernel on distributions by, e.g.,

$$k(\mathbb{P},\mathbb{Q}) = \expigg(-rac{1}{2\sigma^2}\mathrm{MMD}^2(\mathbb{P},\mathbb{Q})igg)$$

• Some pointers: [Muandet+ NeurIPS-12] [Sutherland 2016] [Szabó+ JMLR-16]



Bayesian distribution regression: incorporate $\mu_{\mathbb{P}}$ uncertainty

ightarrow 35



Bayesian distribution regression: incorporate $\mu_{\mathbb{P}}$ uncertainty



IMDb database [Rothe+ 2015]: 400k images of 20k celebrities

Bayesian distribution regression: incorporate $\mu_{\mathbb{P}}$ uncertainty



IMDb database [Rothe+ 2015]: 400k images of 20k celebrities



Bayesian distribution regression: incorporate $\mu_{\mathbb{P}}$ uncertainty



IMDb database [Rothe+ 2015]: 400k images of 20k celebrities



Bayesian distribution regression: incorporate $\mu_{\mathbb{P}}$ uncertainty



• $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g

- $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

 $\mathbb{E}[f(X)] \mathbb{E}[g(Y)]$

- $X \perp\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

 $\mathbb{E}[f(X)] \mathbb{E}[g(Y)] = \langle f, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g \rangle_{\mathcal{H}_y}$

- $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

 $\mathbb{E}[f(X)] \mathbb{E}[g(Y)] = \langle f, \mu_{\mathbb{P}}
angle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g
angle_{\mathcal{H}_y} \ = \langle f, (\mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}})g
angle_{\mathcal{H}_x}$

- $X \perp\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

$$egin{aligned} \mathbb{E}[f(X)] \, \mathbb{E}[g(Y)] &= \langle f, \mu_{\mathbb{P}}
angle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g
angle_{\mathcal{H}_y} \ &= \langle f, (\mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}})g
angle_{\mathcal{H}_x} \ \mathbb{E}[f(X)g(Y)] \end{aligned}$$

- $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

$$\begin{split} \mathbb{E}[f(X)] \, \mathbb{E}[g(Y)] &= \langle f, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g \rangle_{\mathcal{H}_y} \\ &= \langle f, (\mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}}) g \rangle_{\mathcal{H}_x} \\ \mathbb{E}[f(X)g(Y)] &= \mathbb{E}[\langle f, k_x(X, \cdot) \rangle_{\mathcal{H}_x} \langle k_y(Y, \cdot), g \rangle_{\mathcal{H}_y}] \end{split}$$

- $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

 $egin{aligned} \mathbb{E}[f(X)] \ \mathbb{E}[g(Y)] &= \langle f, \mu_{\mathbb{P}}
angle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g
angle_{\mathcal{H}_y} \ &= \langle f, (\mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}})g
angle_{\mathcal{H}_x} \ \mathbb{E}[f(X)g(Y)] &= \mathbb{E}[\langle f, k_x(X, \cdot)
angle_{\mathcal{H}_x} \langle k_y(Y, \cdot), g
angle_{\mathcal{H}_y}] \ &= \langle f, \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \ g
angle_{\mathcal{H}_x} \end{aligned}$
Independence

- $X \perp\!\!\!\!\perp Y$ iff $\operatorname{Cov}(f(X),g(Y)) = 0$ for all measurable f, g
- Let's implement for RKHS functions $f\in \mathcal{H}_x$, $g\in \mathcal{H}_y$:

 $egin{aligned} \mathbb{E}[f(X)] & \mathbb{E}[g(Y)] = \langle f, \mu_{\mathbb{P}}
angle_{\mathcal{H}_x} \langle \mu_{\mathbb{Q}}, g
angle_{\mathcal{H}_y} \ &= \langle f, (\mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}})g
angle_{\mathcal{H}_x} \ & \mathbb{E}[f(X)g(Y)] = \mathbb{E}[\langle f, k_x(X, \cdot)
angle_{\mathcal{H}_x} \langle k_y(Y, \cdot), g
angle_{\mathcal{H}_y}] \ &= \langle f, \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \, g
angle_{\mathcal{H}_x} \ & \mathrm{Cov}(f(X), g(Y)) = \langle f, C_{XY}g
angle_{\mathcal{H}_x} \end{aligned}$

where $C_{XY}: \mathcal{H}_y
ightarrow \mathcal{H}_x$ is

 $\mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)]-\mathbb{E}[k_x(X,\cdot)]\otimes \mathbb{E}[k_y(Y,\cdot)]$

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}}$

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_{\mathbb{P}} \otimes \mu_{\mathbb{Q}}$
- If $X \perp\!\!\!\perp Y$, then $C_{XY} = 0$

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_\mathbb{P} \otimes \mu_\mathbb{Q}$
- If $X \perp\!\!\!\perp Y$, then $C_{XY} = 0$
- If $C_{XY}=0$, $\operatorname{Cov}(f(X),g(Y))=0 \quad orall f \in \mathcal{H}_x, g \in \mathcal{H}_y$

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_\mathbb{P} \otimes \mu_\mathbb{Q}$
- If $X \perp\!\!\!\perp Y$, then $C_{XY} = 0$
- If $C_{XY}=0$, $\operatorname{Cov}(f(X),g(Y))=0 \quad orall f \in \mathcal{H}_x, g \in \mathcal{H}_y$
- If k_x , k_y are characteristic:
 - $C_{XY} = 0$ implies $X \perp \!\!\!\perp Y$ [Szabó/Sriperumbudur JMLR-18]

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_\mathbb{P} \otimes \mu_\mathbb{Q}$
- If $X \perp\!\!\!\perp Y$, then $C_{XY} = 0$
- If $C_{XY}=0$, $\operatorname{Cov}(f(X),g(Y))=0 \quad orall f \in \mathcal{H}_x, g \in \mathcal{H}_y$
- If k_x , k_y are characteristic:
 - $C_{XY} = 0$ implies $X \perp \!\!\!\perp Y$ [Szabó/Sriperumbudur JMLR-18]
 - $X \perp\!\!\!\perp Y$ iff $C_{XY} = 0$

- $\operatorname{Cov}(f(X), g(Y)) = \langle f, C_{XY}g \rangle_{\mathcal{H}_x}$
- $C_{XY} = \mathbb{E}[k_x(X, \cdot) \otimes k_y(Y, \cdot)] \mu_\mathbb{P} \otimes \mu_\mathbb{Q}$
- If $X \perp\!\!\!\perp Y$, then $C_{XY} = 0$
- If $C_{XY}=0$, $\operatorname{Cov}(f(X),g(Y))=0 \quad orall f \in \mathcal{H}_x, g \in \mathcal{H}_y$
- If k_x , k_y are characteristic:
 - $C_{XY} = 0$ implies $X \perp \!\!\!\perp Y$ [Szabó/Sriperumbudur JMLR-18]
 - $X \perp\!\!\!\perp Y$ iff $C_{XY} = 0$
 - $X \perp Y$ iff $0 = \|C_{XY}\|_{\mathrm{HS}}^2$ (sum squared singular values) \circ HSIC: "Hilbert-Schmidt Independence Criterion"

$egin{aligned} C_{XY} &= \mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)] - \mu_\mathbb{P}\otimes \mu_\mathbb{Q} \ \|C_{XY}\|_{\mathrm{HS}}^2 &= \|\mu_{\mathbb{P}_{XY}} - \mu_\mathbb{P}\otimes \mu_\mathbb{Q}\|_{\mathcal{H}_x\otimes\mathcal{H}_y}^2 \end{aligned}$

$egin{aligned} C_{XY} &= \mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)] - \mu_\mathbb{P}\otimes \mu_\mathbb{Q} \ \|C_{XY}\|_{\mathrm{HS}}^2 &= \|\mu_{\mathbb{P}_{XY}} - \mu_\mathbb{P}\otimes \mu_\mathbb{Q}\|_{\mathcal{H}_x\otimes\mathcal{H}_y}^2 \ &= \mathrm{MMD}(\mathbb{P}_{XY},\mathbb{P} imes\mathbb{Q})^2 \end{aligned}$

$$egin{aligned} C_{XY} &= \mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)]-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\ &\|C_{XY}\|^2_{\mathrm{HS}} = \|\mu_{\mathbb{P}_{XY}}-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\|^2_{\mathcal{H}_x\otimes\mathcal{H}_y}\ &= \mathrm{MMD}(\mathbb{P}_{XY},\mathbb{P} imes\mathbb{Q})^2\ &= \mathbb{E}[k_x(X,X')k_y(Y,Y')]\ &-2\,\mathbb{E}[k_x(X,X')k_x(Y,Y'')]\ &+ \mathbb{E}[k_x(X,X')]\,\mathbb{E}[k_y(Y,Y')] \end{aligned}$$

$$egin{aligned} C_{XY} &= \mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)]-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\ \|C_{XY}\|^2_{\mathrm{HS}} &= \|\mu_{\mathbb{P}_{XY}}-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\|^2_{\mathcal{H}_x\otimes\mathcal{H}_y}\ &= \mathrm{MMD}(\mathbb{P}_{XY},\mathbb{P} imes\mathbb{Q})^2\ &= \mathbb{E}[k_x(X,X')k_y(Y,Y')]\ &-2\,\mathbb{E}[k_x(X,X')k_x(Y,Y'')]\ &+ \mathbb{E}[k_x(X,X')]\,\mathbb{E}[k_y(Y,Y')] \end{aligned}$$

• Linear case: C_{XY} is cross-covariance matrix, HSIC is squared Frobenius norm

$$egin{aligned} C_{XY} &= \mathbb{E}[k_x(X,\cdot)\otimes k_y(Y,\cdot)]-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\ \|C_{XY}\|^2_{\mathrm{HS}} &= \|\mu_{\mathbb{P}_{XY}}-\mu_\mathbb{P}\otimes \mu_\mathbb{Q}\|^2_{\mathcal{H}_x\otimes\mathcal{H}_y}\ &= \mathrm{MMD}(\mathbb{P}_{XY},\mathbb{P} imes\mathbb{Q})^2\ &= \mathbb{E}[k_x(X,X')k_y(Y,Y')]\ &-2\,\mathbb{E}[k_x(X,X')k_x(Y,Y'')]\ &+ \mathbb{E}[k_x(X,X')]\,\mathbb{E}[k_y(Y,Y')] \end{aligned}$$

- Linear case: C_{XY} is cross-covariance matrix, HSIC is squared Frobenius norm
- Default estimator (biased, but simple): $\operatorname{Tr}(HK_XHK_Y)$ where $H = I - \mathbf{11}^T$

HSIC applications

- Independence testing [Gretton+ NeurIPS-07]
- Clustering [Song+ ICML-07]
- Feature selection [Song+ JMLR-12]
- Self-supervised learning [Li+ NeurIPS-21]
- :
- Broadly: easier-to-estimate, sometimes-nicer version of mutual information

Example: SSL-HSIC [Li+ NeurIPS-21]



- Maximizes dependence between image features $m{f}$ and its identity on a minibatch
- Using a learned deep kernel based on $oldsymbol{g}$

Recap

- Mean embedding $\mu_{\mathbb{P}} = \mathbb{E} \, k(X, \cdot)$
- $MMD(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} \mu_{\mathbb{Q}}\|_{\mathcal{H}}$ is 0 iff $\mathbb{P} = \mathbb{Q}$ (for characteristic kernels)
- $\operatorname{HSIC}(X, Y) = \|C_{XY}\|_{\operatorname{HS}} = \operatorname{MMD}(\mathbb{P}_{XY}, \mathbb{P} \times \mathbb{Q})^2$ is 0 iff $X \perp \!\!\!\perp Y$ (for characteristic k_x , k_y)
- After break: last interactive session exploring testing
- More details:
 - Close connections to Gaussian processes [Kanagawa+ 'GPs and Kernel Methods' 2018]
 - Mean embeddings: survey [Muandet+ 'Kernel Mean Embedding of Distributions']