Better deep learning (sometimes) by learning kernel mean embeddings

Danica J. Sutherland (she/her)

University of British Columbia (UBC) / Alberta Machine Intelligence Institute (Amii)

slides (but not the talk) about four related projects:



LIKE22 - 12 Jan 2022

- Deep learning: models usually of form $f(x) = w^{\mathsf{T}} \phi_{\psi}(x)$
 - With a learned $\phi_\psi(x):\mathcal{X} o\mathbb{R}^D$

- Deep learning: models usually of form $f(x) = w^{\mathsf{T}} \phi_{\psi}(x)$ With a *learned* $\phi_{\psi}(x) : \mathcal{X} \to \mathbb{R}^D$
- 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: models usually of form $f(x) = w^{\mathsf{T}} \phi_{\psi}(x)$ With a *learned* $\phi_{\psi}(x) : \mathcal{X} \to \mathbb{R}^D$
- 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: models usually of form $f(x) = w^{\mathsf{T}} \phi_{\psi}(x)$ With a *learned* $\phi_{\psi}(x) : \mathcal{X} \to \mathbb{R}^D$
- 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
- Could train a classifier by:
 - Let $ilde{L}(\psi) = L(f_\psi^*)$, loss of the best $f_\psi^* \in H_\psi$

Learn ψ by following $abla_{\psi} ilde{L}(\psi) =
abla_{\psi} L(f_{\psi}^*)$

- Deep learning: models usually of form $f(x) = w^{\mathsf{T}} \phi_{\psi}(x)$ With a *learned* $\phi_{\psi}(x) : \mathcal{X} \to \mathbb{R}^D$
- 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
- Could train a classifier by:

- Let $ilde{L}(\psi) = L(f_\psi^*)$, loss of the best $f_\psi^* \in H_\psi$

Learn ψ by following $abla_{\psi} ilde{L}(\psi) =
abla_{\psi} L(f_{\psi}^*)$

• Generalize to a **deep kernel**:

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

- Take $\phi_\psi(x) \in \mathbb{R}$ as output of *last* layer
- $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)$
- Final function in \mathcal{H}_ψ will be $a\phi_\psi(x)$

- Take $\phi_\psi(x) \in \mathbb{R}$ as output of *last* layer
- $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)+1$
- Final function in \mathcal{H}_ψ will be $a\phi_\psi(x)+b$

- Take $\phi_\psi(x) \in \mathbb{R}$ as output of *last* layer
- $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)+1$
- Final function in \mathcal{H}_ψ will be $a\phi_\psi(x)+b$
- With logistic loss: this is Platt scaling

- Take $\phi_\psi(x) \in \mathbb{R}$ as output of *last* layer
- $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)+1$
- Final function in \mathcal{H}_ψ will be $a\phi_\psi(x)+b$
- With logistic loss: this is Platt scaling

On Calibration of Modern Neural Networks

Chuan Guo^{*1} Geoff Pleiss^{*1} Yu Sun^{*1} Kilian Q. Weinberger¹

• 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 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 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 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

• In "normal" classification: slightly richer function space 🖓

- In "normal" classification: slightly richer function space 🦓
- Meta-learning: common ϕ_ψ , constantly varying f^*_ψ

META-LEARNING WITH DIFFERENTIABLE CLOSED-FORM SOLVERS

Luca Bertinetto FiveAI & University of Oxford luca@robots.ox.ac.uk

Philip H.S. Torr
FiveAI & University of Oxford
philip.torr@eng.ox.ac.uk

João Henriques University of Oxford joao@robots.ox.ac.uk

Andrea Vedaldi University of Oxford vedaldi@robots.ox.ac.uk

Meta-Learning with Differentiable Convex Optimization

Kwonjoon Lee2SubhransuMajiAvinashRavichandranStefano Soatto1Amazon Web Services2UC San Diego3UMass Amherst4UCLAkwl042@ucsd.edu{smmaji,ravinash,soattos}@amazon.com

- In "normal" classification: slightly richer function space
- Meta-learning: common ϕ_ψ , constantly varying f_ψ^*

- In "normal" classification: slightly richer function space
- Meta-learning: common ϕ_ψ , constantly varying f_ψ^*
- Two-sample testing
 - Simple form of f_ψ^* for cheap permutation testing

- In "normal" classification: slightly richer function space
- Meta-learning: common ϕ_ψ , constantly varying f_ψ^*
- Two-sample testing
 - Simple form of f_ψ^* for cheap permutation testing
- Self-supervised learning
 - Better understanding of what's really going on, at least

- In "normal" classification: slightly richer function space
- Meta-learning: common ϕ_ψ , constantly varying f_ψ^*
- Two-sample testing
 - Simple form of f_{ψ}^{*} for cheap permutation testing
- Self-supervised learning
 - Better understanding of what's really going on, at least
- Generative modeling with MMD GANs
 - Better gradient for generator to follow (?)

- In "normal" classification: slightly richer function space 🖓
- Meta-learning: common ϕ_ψ , constantly varying f_ψ^*
- Two-sample testing
 - Simple form of f_ψ^* for cheap permutation testing
- Self-supervised learning
 - Better understanding of what's really going on, at least
- Generative modeling with MMD GANs
 Better gradient for generator to follow (?)
- Score matching in exponential families (density estimation)
 - Optimize regularization weights, better gradient (?)

$\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} [f(X)] - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} [f(Y)]$

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

 $egin{aligned} \mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}\left[f(X)
ight] - \mathbb{E}\left[f(Y)
ight] \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}\left[\langle f, arphi(X)
angle_{\mathcal{H}}
ight] - \mathbb{E}\left[\langle f, arphi(Y)
angle_{\mathcal{H}}
ight] \ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \left\langle f, \mathbb{E}\left[arphi(X)
ight] - \mathbb{E}\left[arphi(X)
ight] - \mathbb{E}\left[arphi(Y)
ight]
angle_{\mathcal{H}} \end{aligned}$

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

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

MMD as feature matching

$$\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \left\| \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[arphi(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[arphi(Y)]
ight\|_{\mathcal{H}}$$

. .

• $arphi: X o \mathcal{H}$ is the *feature map* for $k(x,y) = \langle arphi(x), arphi(y)
angle$

MMD as feature matching

$$\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \left\| \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[arphi(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[arphi(Y)]
ight\|_{\mathcal{H}}$$

- $arphi: X o \mathcal{H}$ is the *feature map* for $k(x,y) = \langle arphi(x), arphi(y)
 angle$
- If $k(x,y) = x^{\mathsf{T}}y$, $\varphi(x) = x$, then the MMD is distance between means

MMD as feature matching

$$\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \left\| \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[arphi(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[arphi(Y)]
ight\|_{\mathcal{H}}$$

- $arphi: X o \mathcal{H}$ is the *feature map* for $k(x,y) = \langle arphi(x), arphi(y)
 angle$
- If $k(x,y) = x^{\mathsf{T}}y$, $\varphi(x) = x$, then the MMD is distance between means
- Many kernels: infinite-dimensional ${\cal H}$

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

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

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

$$egin{aligned} \mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) &= \mathop{\mathbb{E}}_{X,X'\sim\mathbb{P}}[k(X,X')] + \mathop{\mathbb{E}}_{Y,Y'\sim\mathbb{Q}}[k(Y,Y')] - 2 \mathop{\mathbb{E}}_{\substack{X\sim\mathbb{P}}{Y\sim\mathbb{Q}}}[k(X,Y)] \ & \frown \ 2 \end{aligned}$$

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

K_{XX}



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

 $\widetilde{\mathrm{MMD}}_{k}(X,Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \mathrm{mean}(K_{XY})$

 K_{XX}

$$K_{YY}$$

1.0	0.2	0.6	(<u> </u>	1.0	0.8	0.7
0.2	1.0	0.5		0.8	1.0	0.6
0.6	0.5	1.0	- 2000 - 2000).	0.7	0.6	1.0

$$egin{aligned} \mathrm{MMD}_k^2(\mathbb{P},\mathbb{Q}) &= \mathop{\mathbb{E}}_{X,X'\sim\mathbb{P}}[k(X,X')] + \mathop{\mathbb{E}}_{Y,Y'\sim\mathbb{Q}}[k(Y,Y')] - 2 \mathop{\mathbb{E}}_{\substack{X\sim\mathbb{P}\\Y\sim\mathbb{Q}}}[k(X,Y)] \ & \widehat{\mathrm{MMD}}_k^2(X,Y) = \mathrm{mean}(K_{XX}) + \mathrm{mean}(K_{YY}) - 2 \mathop{\mathrm{mean}}(K_{XY}) \end{aligned}$$



1.0	0.2	0.6		1.0	0.8	0.7		0.3	0.1	0.2
0.2	1.0	0.5		0.8	1.0	0.6	· <u>(</u>),	0.2	0.3	0.3
0.6	0.5	1.0	(Case _ ass)	0.7	0.6	1.0	(<u>Cano, _ mar</u>))	0.2	0.1	0.4

I: Two-sample testing

• Given samples from two unknown distributions

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

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

I: Two-sample testing

• 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



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

• 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?

• 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

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

- 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

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

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

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

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

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

- 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?
- 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} ?
- Independence testing: is P(X, Y) = P(X)P(Y)?

• 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:

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

- Reject H_0 if test statistic $\hat{T}(X,Y)>c_lpha$











 $\begin{array}{l} \text{Permutation testing to find } c_{\alpha} \\ & \text{Need } \Pr_{H_0} \left(T(X,Y) > c_{\alpha} \right) \leq \alpha \\ X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5 \qquad Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \\ c_{\alpha} : 1 - \alpha \text{th quantile of } \left\{ \begin{array}{c} & & \end{array} \right\} \end{array}$







 $\begin{array}{l} \text{Permutation testing to find } c_{\alpha} \\ \text{Need } \Pr_{H_0} \left(T(X,Y) > c_{\alpha} \right) \leq \alpha \\ X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5 \quad Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5 \\ c_{\alpha} : 1 - \alpha \text{th quantile of } \left\{ \hat{T}(\tilde{X}_1,\tilde{Y}_1), \ \hat{T}(\tilde{X}_2,\tilde{Y}_2), \ \cdots \right\} \end{array}$

- If k is characteristic, $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
- Efficient permutation testing for $\widehat{\mathrm{MMD}}(X,Y)$

- If k is characteristic, $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
- Efficient permutation testing for $\widehat{\mathrm{MMD}}(X,Y)$
 - $H_0: \widehat{\mathrm{nMMD}}^2$ converges in distribution
 - $H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 \mathrm{MMD}^2)$ asymptotically normal

- If k is characteristic, $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
- Efficient permutation testing for $\widehat{\mathrm{MMD}}(X,Y)$
 - $H_0: n\widehat{\mathrm{MMD}}^2$ converges in distribution
 - $H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 \mathrm{MMD}^2)$ asymptotically normal
- Any characteristic kernel gives consistent test

- If k is characteristic, $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
- Efficient permutation testing for $\widehat{\mathrm{MMD}}(X,Y)$
 - $H_0: n\widehat{\mathrm{MMD}}^2$ converges in distribution
 - $H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 \mathrm{MMD}^2)$ asymptotically normal
- Any characteristic kernel gives consistent test...eventually

- If k is characteristic, $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$ iff $\mathbb{P}=\mathbb{Q}$
- Efficient permutation testing for $\widehat{\mathrm{MMD}}(X,Y)$
 - $H_0: n\widehat{\mathrm{MMD}}^2$ converges in distribution
 - $H_1: \sqrt{n}(\widehat{\mathrm{MMD}}^2 \mathrm{MMD}^2)$ 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 $\widehat{\mathrm{MMD}}(X,Y)=\left|\hat{T}(X,Y)-rac{1}{2}
 ight|$

Optimizing test power

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

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

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

Optimizing test power

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

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

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

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

Optimizing test power

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

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

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

- Pick k to maximize an estimate of $\mathrm{MMD}^2 \, / \sigma_{H_1}$
- Can show uniform $\mathcal{O}_P(n^{-rac{1}{3}})$ convergence of estimator

Blobs dataset



Blobs kernels



Blobs results



CIFAR-10 vs CIFAR-10.1



Train on 1 000, test on 1 031, repeat 10 times. Rejection rates:

ME	SCF	C2ST	MMD-O	MMD-D
0.588	0.171	0.452	0.316	0.744

Ablation vs classifier-based tests

	Cross-entropy			Max power		
Dataset	Sign	Lin	Ours	Sign	Lin	Ours
Blobs	0.84	0.94	0.90	_	0.95	0.99
High- d Gauss. mix.	0.47	0.59	0.29	_	0.64	0.66
Higgs	0.26	0.40	0.35	_	0.30	0.40
MNIST vs GAN	0.65	0.71	0.80	_	0.94	1.00

But...

• What if you don't have much data for your testing problem?
- What if you don't have much data for your testing problem?
- Need enough data to pick a good kernel

- What if you don't have much data for your testing problem?
- Need enough data to pick a good kernel
- Also need enough test data to actually detect the difference

- What if you don't have much data for your testing problem?
- Need enough data to pick a good kernel
- Also need enough test data to actually detect the difference
- Best split depends on best kernel's quality / how hard to find

- What if you don't have much data for your testing problem?
- Need enough data to pick a good kernel
- Also need enough test data to actually detect the difference
- Best split depends on best kernel's quality / how hard to find
 - Don't know that ahead of time; can't try more than one

Meta-testing

• One idea: what if we have *related* problems?

Meta-testing

- One idea: what if we have *related* problems?
- Similar setup to meta-learning:



Meta-testing for CIFAR-10 vs CIFAR-10.1

- CIFAR-10 has 60,000 images, but CIFAR-10.1 only has 2,031
- Where do we get related data from?

Meta-testing for CIFAR-10 vs CIFAR-10.1

- CIFAR-10 has 60,000 images, but CIFAR-10.1 only has 2,031
- Where do we get related data from?
- One option: set up tasks to distinguish classes of CIFAR-10 (airplane vs automobile, airplane vs bird, ...)

One approach (MAML-like)



One approach (MAML-like)



This works, but not as well as we'd hoped... Initialization might work okay on everything, not really adapt

Another approach: Meta-MKL



Inspired by classic multiple kernel learning

Only need to learn linear combination β_i on test task: much easier

Testing Samples

Meta-Samples

Theoretical analysis for Meta-MKL

- Same big-O dependence on test task size 😐
- But multiplier is *much* better: based on number of meta-training tasks, not on network size

Theoretical analysis for Meta-MKL

- Same big-O dependence on test task size 😐
- But multiplier is *much* better: based on number of meta-training tasks, not on network size
- (Analysis assumes meta-tasks are "related" enough)

Results on CIFAR-10.1

Methods	$m_{tr} = 100$			$m_{tr} = 200$		
	$m_{te} = 200$	$m_{te} = 500$	$m_{te} = 900$	$m_{te} = 200$	$m_{te} = 500$	$m_{te} = 900$
ME	$0.084 \scriptstyle \pm 0.009$	$0.096 \scriptstyle \pm 0.016$	$0.160{\scriptstyle \pm 0.035}$	$0.104 \scriptstyle \pm 0.013$	$0.202 \scriptstyle \pm 0.020$	$0.326 \scriptscriptstyle \pm 0.039$
SCF	$0.047 \scriptstyle \pm 0.013$	$0.037 \scriptstyle \pm 0.011$	$0.047 \scriptstyle \pm 0.015$	$0.026 \scriptstyle \pm 0.009$	$0.018 \scriptstyle \pm 0.006$	$0.026 \scriptstyle \pm 0.012$
C2ST-S	$0.059 \scriptstyle \pm 0.009$	$0.062 \scriptstyle \pm 0.007$	$0.059 \scriptstyle \pm 0.007$	$0.052 \scriptstyle \pm 0.011$	$0.054 \scriptscriptstyle \pm 0.011$	$0.057 \scriptstyle \pm 0.008$
C2ST-L	$0.064 \scriptstyle \pm 0.009$	$0.064 \scriptstyle \pm 0.006$	$0.063 \scriptstyle \pm 0.007$	$0.075 \scriptstyle \pm 0.014$	$0.066 \scriptstyle \pm 0.011$	$0.067 \scriptstyle \pm 0.008$
MMD-O	$0.091 \scriptstyle \pm 0.011$	$0.141 \scriptstyle \pm 0.009$	$0.279 \scriptstyle \pm 0.018$	$0.084 \scriptstyle \pm 0.007$	$0.160 \scriptstyle \pm 0.011$	$0.319 \scriptstyle \pm 0.020$
MMD-D	$0.104 \scriptstyle \pm 0.007$	$0.222{\scriptstyle \pm 0.020}$	$0.418 \scriptstyle \pm 0.046$	$0.117 \scriptstyle \pm 0.013$	$0.226 \scriptstyle \pm 0.021$	$0.444{\scriptstyle \pm 0.037}$
AGT-KL	$0.170_{\pm 0.032}$	0.457	$0.765 \scriptstyle \pm 0.045$	0.152	0.463 ± 0.060	$0.778 \scriptstyle \pm 0.050$
Meta-KL	$0.245{\scriptstyle \pm 0.010}$	$0.671 \scriptstyle \pm 0.026$	$0.959 \scriptstyle \pm 0.013$	$0.226 \scriptstyle \pm 0.015$	$0.668 \scriptstyle \pm 0.032$	$0.972 \scriptstyle \pm 0.006$
Meta-MKL	$0.277 \scriptscriptstyle \pm 0.016$	$0.728 \scriptstyle \pm 0.020$	$0.973 \scriptscriptstyle \pm 0.008$	$0.255 \scriptstyle \pm 0.020$	$0.724 \scriptscriptstyle \pm 0.026$	$0.993 \scriptstyle \pm 0.003$

• When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?

- When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?
 - Methods so far: some mostly for low-d
 - Some look at points with large critic function

- When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?
 - Methods so far: some mostly for low-d
 - Some look at points with large critic function



- When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?
 - Methods so far: some mostly for low-d
 - Some look at points with large critic function



- Finding kernels / features that *can't* do certain things
 - distinguish by emotion, but can't distinguish by skin color

- When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?
 - Methods so far: some mostly for low-d
 - Some look at points with large critic function



- Finding kernels / features that *can't* do certain things
 - distinguish by emotion, but can't distinguish by skin color
- Avoid the need for data splitting (selective inference)

- When $\mathbb{P} \neq \mathbb{Q}$, can we tell *how* they're different?
 - Methods so far: some mostly for low-d
 - Some look at points with large critic function



- Finding kernels / features that *can't* do certain things
 - distinguish by emotion, but can't distinguish by skin color
- Avoid the need for data splitting (selective inference)
 - Kübler+ NeurIPS-20 gave one method
 - only for multiple kernel learning
 - only with data-inefficient (streaming) estimator

II: Self-supervised learning

Given a bunch of unlabeled samples X, want to find "good" features Z=f(X)

(e.g. so that a linear classifier on Z works with few samples)

II: Self-supervised learning

Given a bunch of unlabeled samples X, want to find "good" features Z=f(X)

(e.g. so that a linear classifier on Z works with few samples)

One common approach: contrastive learning



InfoNCE [van den Oord+ 2018, Poole+ ICML-19]

$$\mathbb{E}_{Z_1}iggl[\mathbb{E}_{Z_2\sim ext{pos}}[k(Z_1,Z_2)]] + \log \mathbb{E}_{Z_2}[\exp(k(Z_1,Z_2))iggr] \ \leq \operatorname{MI}(Z_1,Z_2)$$



Variants: CPC, SimCLR, MoCo, SwAV, ...

Mutual information isn't why SSL works!

InfoNCE approximates MI between "positive" views But MI is invariant to transformations important to SSL!



$egin{aligned} \mathrm{HSIC}(X,Y) &= \left\| \mathbb{E}[\phi(X)\otimes\phi(Y)] - \mathbb{E}[\phi(X)]\otimes\mathbb{E}[\phi(Y)] ight\|_{HS}^2 \ &= \mathrm{MMD}^2(\mathbb{P}_{XY},\mathbb{P}_X\otimes\mathbb{P}_Y) \end{aligned}$

$$\begin{split} \mathrm{HSIC}(X,Y) &= \|\mathbb{E}[\phi(X)\otimes\phi(Y)] - \mathbb{E}[\phi(X)]\otimes\mathbb{E}[\phi(Y)]\|_{HS}^2 \\ &= \mathrm{MMD}^2(\mathbb{P}_{XY},\mathbb{P}_X\otimes\mathbb{P}_Y) \\ &\leq C \ \mathrm{MI}(X,Y) = (C \ \mathrm{denom}\,\mathrm{d$$

 $\leq C_k \operatorname{MI}(X,Y) \qquad (C_k ext{ depends only on } \|k\|_\infty)$

$egin{aligned} \mathrm{HSIC}(X,Y) &= ig\|\mathbb{E}[\phi(X)\otimes\phi(Y)] - \mathbb{E}[\phi(X)]\otimes\mathbb{E}[\phi(Y)]ig\|_{HS}^2 \ &= \mathrm{MMD}^2(\mathbb{P}_{XY},\mathbb{P}_X\otimes\mathbb{P}_Y) \ &\leq C_k\,\mathrm{MI}(X,Y) \qquad (C_k ext{ depends only on }\|k\|_\infty) \end{aligned}$

With a linear kernel: $\operatorname{HSIC} = \left\| \mathbb{E}[XY^{\mathsf{T}}] - \mathbb{E}[X] \mathbb{E}[Y]^{\mathsf{T}} \right\|_{F}^{2}$

$egin{aligned} \mathrm{HSIC}(X,Y) &= \|\mathbb{E}[\phi(X)\otimes\phi(Y)] - \mathbb{E}[\phi(X)]\otimes\mathbb{E}[\phi(Y)]\|_{HS}^2 \ &= \mathrm{MMD}^2(\mathbb{P}_{XY},\mathbb{P}_X\otimes\mathbb{P}_Y) \ &\leq C_k\,\mathrm{MI}(X,Y) \qquad (C_k ext{ depends only on }\|k\|_\infty) \end{aligned}$

With a linear kernel: $\operatorname{HSIC} = \left\| \mathbb{E}[XY^{\mathsf{T}}] - \mathbb{E}[X] \mathbb{E}[Y]^{\mathsf{T}} \right\|_{F}^{2}$

Estimator based on kernel matrices:

$$\widehat{\mathrm{HSIC}} = rac{1}{(n-1)^2} \mathrm{tr}(KHLH)$$

- $oldsymbol{H}$ is the centering matrix
- old K is the kernel matrix on old X
- L is the kernel matrix on Y

SSL-HSIC

 $\mathcal{L}_{ ext{SSL-HSIC}} = -\operatorname{HSIC}(Z, Y) + \gamma \sqrt{\operatorname{HSIC}(Z, Z)}$

(Y is just an indicator of which source image it came from)



Target representation Z' is output of $f_ heta(X)$

HSIC uses learned kernel $k(Z_1',Z_2') = \mathrm{IMQ}(g(Z_1'),g(Z_2'))$

$\mathcal{L}_{ ext{InfoNCE}}(heta) = - \mathop{\mathbb{E}}_{(Z_1,Z_2)\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\log\mathop{\mathbb{E}}_{Z_2}[\exp k(Z_1,Z_2)]$

$$egin{aligned} \mathcal{L}_{ ext{InfoNCE}}(heta) &= - \mathop{\mathbb{E}}_{(Z_1,Z_2)\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\log\mathop{\mathbb{E}}_{Z_2}[\exp{k(Z_1,Z_2)}] \ &pprox &= -\mathop{\mathbb{E}}_{z_1,z_2\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\mathop{\mathbb{E}}_{Z_2}[k(Z_1,Z_2)] + rac{1}{2}\mathop{\mathbb{E}}_{Z_1}\left[\mathop{ ext{Var}}_{Z_2}[k(Z_1,Z_2)]
ight] \ &pprox &= -\mathop{ ext{HSIC}}_{Z_2}[k(Z_1,Z_2)] + \mathop{ ext{HSIC}}_{Z_1}\left[k(Z_1,Z_2)\right] + \mathop{ ext{HSIC}}_{Z_2}\left[k(Z_1,Z_2)\right] + \mathop{ ext{HSI$$

$$egin{aligned} \mathcal{L}_{ ext{InfoNCE}}(heta) &= - \mathop{\mathbb{E}}_{(Z_1,Z_2)\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\log\mathop{\mathbb{E}}_{Z_2}[\exp{k(Z_1,Z_2)}] \ &pprox &= -\mathop{\mathbb{E}}_{z_1,z_2\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\mathop{\mathbb{E}}_{Z_2}[k(Z_1,Z_2)] + rac{1}{2}\mathop{\mathbb{E}}_{Z_1}\left[\mathop{ ext{Var}}_{Z_2}[k(Z_1,Z_2)]
ight] \ &pprox &= -\mathop{ ext{HSIC}}_{Z_2}[k(Z_1,Z_2)] + \mathop{ ext{E}}_{Z_1}\mathcal{E}_{Z_2}[k(Z_1,Z_2)] + rac{1}{2}\mathop{ ext{E}}_{Z_2}\left[\mathop{ ext{Var}}_{Z_2}[k(Z_1,Z_2)]
ight] \ & ext{variance penalty} \end{aligned}$$

Very similar loss! Just different regularizer

$$egin{aligned} \mathcal{L}_{ ext{InfoNCE}}(heta) &= - \mathop{\mathbb{E}}_{(Z_1,Z_2)\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\log\mathop{\mathbb{E}}_{Z_2}[\exp{k(Z_1,Z_2)}] \ &pprox &= -\mathop{\mathbb{E}}_{z_1,z_2\sim ext{pos}}[k(Z_1,Z_2)] + \mathop{\mathbb{E}}_{Z_1}\mathop{\mathbb{E}}_{Z_2}[k(Z_1,Z_2)] + rac{1}{2}\mathop{\mathbb{E}}_{Z_1}\left[\mathop{ ext{Var}}_{Z_2}[k(Z_1,Z_2)]
ight] \ &pprox &= \mathop{ ext{MSIC}}_{Z_2}(Z,Y) \end{aligned}$$

Very similar loss! Just different regularizer

When variance is small, $-\operatorname{HSIC}(Z, Y) + \gamma \operatorname{HSIC}(Z, Z) \leq \mathcal{L}_{\operatorname{InfoNCE}} + o(\operatorname{variance})$

Clustering interpretation

SSL-HSIC estimates agreement of Z with cluster structure of Y

With linear kernels:

$$-\operatorname{HSIC}(Z,Y) \propto \sum_{i=1}^n \sum_{p=1}^m \left\|Z_i^{(p)} - ar{Z}_i
ight\|^2 - nm$$

where \bar{Z}_i is mean of the m augmentations

Clustering interpretation

SSL-HSIC estimates agreement of Z with cluster structure of Y

With linear kernels:

$$-\operatorname{HSIC}(Z,Y) \propto \sum_{i=1}^n \sum_{p=1}^m \left\| Z_i^{(p)} - ar Z_i
ight\|^2 - nm$$

where \overline{Z}_i is mean of the m augmentations

Resembles BYOL loss with no target network (but still works!)

ImageNet results: linear evaluation


Transfer from ImageNet to classification tasks



III: Training implicit generative models

Given samples from a distribution \mathbb{P} over \mathcal{X} , we want a model that can produce new samples from $\mathbb{Q}_{\theta} \approx \mathbb{P}$



 $X\sim \mathbb{P}$

 $Y \sim \mathbb{Q}_{\theta}$

III: Training implicit generative models

we want



 $\mathbb{P}_{ heta} pprox \mathbb{P}$

thispersondoesnotexist.com

III: Training implicit generative models

Given samples from a distribution \mathbb{P} over \mathcal{X} , we want a model that can produce new samples from $\mathbb{Q}_{\theta} \approx \mathbb{P}$



Generator networks

Fixed distribution of latents: $Z\sim ext{Uniform}\left([-1,1]^{100}
ight)$ Maps through a network: $G_ heta(Z)\sim \mathbb{Q}_ heta$



DCGAN generator [Radford+ ICLR-16]

Generator networks

Fixed distribution of latents: $Z\sim ext{Uniform}\left([-1,1]^{100}
ight)$ Maps through a network: $G_{ heta}(Z)\sim \mathbb{Q}_{ heta}$



DCGAN generator [Radford+ ICLR-16] How to choose θ ?

- GANs [Goodfellow+ NeurIPS-14] minimize discriminator accuracy (like classifier test) between \mathbb{P} and \mathbb{Q}_{θ}
- Problem: if there's a perfect classifier, discontinuous loss, no gradient to improve it [Arjovsky/Bottou ICLR-17]
- Disjoint at init:





- GANs [Goodfellow+ NeurIPS-14] minimize discriminator accuracy (like classifier test) between \mathbb{P} and \mathbb{Q}_{θ}
- Problem: if there's a perfect classifier, discontinuous loss, no gradient to improve it [Arjovsky/Bottou ICLR-17]
- Disjoint at init:



• For usual $G_ heta: \mathbb{R}^{100} o \mathbb{R}^{64 imes 64 imes 3}$, $\mathbb{Q}_ heta$ is supported on a countable union of manifolds with dim ≤ 100

- GANs [Goodfellow+ NeurIPS-14] minimize discriminator accuracy (like classifier test) between \mathbb{P} and \mathbb{Q}_{θ}
- Problem: if there's a perfect classifier, discontinuous loss, no gradient to improve it [Arjovsky/Bottou ICLR-17]
- Disjoint at init:



- For usual $G_ heta: \mathbb{R}^{100} o \mathbb{R}^{64 imes 64 imes 3}$, $\mathbb{Q}_ heta$ is supported on a countable union of manifolds with dim ≤ 100
- "Natural image manifold" usually considered low-dim

- GANs [Goodfellow+ NeurIPS-14] minimize discriminator accuracy (like classifier test) between \mathbb{P} and \mathbb{Q}_{θ}
- Problem: if there's a perfect classifier, discontinuous loss, no gradient to improve it [Arjovsky/Bottou ICLR-17]
- Disjoint at init:



- For usual $G_ heta: \mathbb{R}^{100} o \mathbb{R}^{64 imes 64 imes 3}$, $\mathbb{Q}_ heta$ is supported on a countable union of manifolds with dim ≤ 100
- "Natural image manifold" usually considered low-dim
- Won't align at init, so won't ever align

- Integral probability metrics with "smooth" ${\mathcal F}$ are continuous
- WGAN: ${\mathcal F}$ a set of neural networks satisfying $\|f\|_L \leq 1$
- WGAN-GP: instead penalize $\mathbb{E} \|
 abla_x f(x) \|$ near the data

- Integral probability metrics with "smooth" ${\mathcal F}$ are continuous
- WGAN: ${\mathcal F}$ a set of neural networks satisfying $\|f\|_L \leq 1$
- WGAN-GP: instead penalize $\mathbb{E} \|
 abla_x f(x) \|$ near the data
- Both losses are MMD with $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)$

- Integral probability metrics with "smooth" ${\mathcal F}$ are continuous
- WGAN: ${\mathcal F}$ a set of neural networks satisfying $\|f\|_L \leq 1$
- WGAN-GP: instead penalize $\mathbb{E} \|
 abla_x f(x) \|$ near the data
- Both losses are MMD with $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)$

$$\bullet \quad \min_{\theta} \left[\mathcal{D}^{\Psi}_{\mathrm{MMD}}(\mathbb{P},\mathbb{Q}_{\theta}) = \sup_{\psi \in \Psi} \mathrm{MMD}_{\psi}(\mathbb{P},\mathbb{Q}_{\theta}) \right]$$

- Integral probability metrics with "smooth" ${\mathcal F}$ are continuous
- WGAN: ${\mathcal F}$ a set of neural networks satisfying $\|f\|_L \leq 1$
- WGAN-GP: instead penalize $\mathbb{E} \|
 abla_x f(x) \|$ near the data
- Both losses are MMD with $k_\psi(x,y)=\phi_\psi(x)\phi_\psi(y)$

$$\bullet \quad \min_{\theta} \left[\mathcal{D}^{\Psi}_{\mathrm{MMD}}(\mathbb{P},\mathbb{Q}_{\theta}) = \sup_{\psi \in \Psi} \mathrm{MMD}_{\psi}(\mathbb{P},\mathbb{Q}_{\theta}) \right]$$

• Some kind of constraint on ϕ_ψ is important!































Illustrative problem in \mathbb{R} , DiracGAN [Mescheder+ ICML-18]:

• Just need to stay away from tiny bandwidths ψ • ...deep kernel analogue is hard.

θ

Instead, keep witness function from being too steep

 $k_{w=2}(0, x)$

θ

Illustrative problem in \mathbb{R} , DiracGAN [Mescheder+ ICML-18]:

 $k_{\psi=0.25}(0, x)$

- Just need to stay away from tiny bandwidths ψ
- ...deep kernel analogue is hard.

 $k_{w=2}(0, x)$

- Instead, keep witness function from being too steep
- $\sup_{x} \| \nabla f(x) \|$ would give Wasserstein
 - Nice distance, but hard to estimate

 $k_{\psi=0.01}(0, x)$

Illustrative problem in \mathbb{R} , DiracGAN [Mescheder+ ICML-18]:

 $k_{\psi=0.25}(0, x)$

- Just need to stay away from tiny bandwidths ψ
- ...deep kernel analogue is hard.

 $k_{\psi=2}(0, x)$

- Instead, keep witness function from being too steep
- $\sup_{x} \| \nabla f(x) \|$ would give Wasserstein • Nice distance, but hard to estimate
- Control $\|
 abla f(ilde X)\|$ on average, near the data
 - Gulrajani+ NeurIPS-17 / Roth+ NeurIPS-17 / Mescheder+ ICML-18

 $k_{w=0.01}(0, x)$

- If Ψ gives uniformly Lipschitz critics, $\mathcal{D}^{\Psi}_{\mathrm{MMD}}$ is smooth
- Original MMD-GAN paper [Li+ NeurIPS-17]: box constraint

- If Ψ gives uniformly Lipschitz critics, $\mathcal{D}^{\Psi}_{\mathrm{MMD}}$ is smooth
- Original MMD-GAN paper [Li+ NeurIPS-17]: box constraint
- We [Bińkowski+ ICLR-18] used gradient penalty on critic instead

- If Ψ gives uniformly Lipschitz critics, $\mathcal{D}^{\Psi}_{\mathrm{MMD}}$ is smooth
- Original MMD-GAN paper [Li+ NeurIPS-17]: box constraint
- We [Bińkowski+ ICLR-18] used gradient penalty on critic instead
 - Better in practice, but doesn't fix the Dirac problem...

- If Ψ gives uniformly Lipschitz critics, $\mathcal{D}^{\Psi}_{\mathrm{MMD}}$ is smooth
- Original MMD-GAN paper [Li+ NeurIPS-17]: box constraint
- We [Bińkowski+ ICLR-18] used gradient penalty on critic instead
 - Better in practice, but doesn't fix the Dirac problem...


New distance: Scaled MMD Want to ensure $\mathbb{E}_{ ilde{X}\sim\mathbb{S}}[\| abla f(ilde{X})\|^2]\leq 1$

Want to ensure $\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] \leq 1$

Can solve with $\langle \partial_i \phi(x), f
angle_{\mathcal{H}} = \partial_i f(x)$...but too expensive!

Want to ensure
$$\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] \leq 1$$

Can solve with $\langle \partial_i \phi(x), f
angle_{\mathcal{H}} = \partial_i f(x)$...but too expensive!

Guaranteed if
$$\|f\|_{\mathcal{H}} \leq \sigma_{\mathbb{S},k,\lambda}$$

 $\sigma_{\mathbb{S},k,\lambda} := \left(\lambda + \mathop{\mathbb{E}}_{\tilde{X}\sim\mathbb{S}}\left[k(\tilde{X},\tilde{X}) + [\nabla_1 \cdot \nabla_2 k](\tilde{X},\tilde{X})\right]\right)^{-\frac{1}{2}}$

Want to ensure
$$\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] \leq 1$$

Can solve with $\langle \partial_i \phi(x), f
angle_{\mathcal{H}} = \partial_i f(x)$...but too expensive!

Guaranteed if
$$\|f\|_{\mathcal{H}} \leq \sigma_{\mathbb{S},k,\lambda}$$

 $\sigma_{\mathbb{S},k,\lambda} := \left(\lambda + \mathop{\mathbb{E}}_{\tilde{X}\sim\mathbb{S}}\left[k(\tilde{X},\tilde{X}) + [\nabla_1 \cdot \nabla_2 k](\tilde{X},\tilde{X})\right]\right)^{-\frac{1}{2}}$

Gives distance $\mathrm{SMMD}_{\mathbb{S},k,\lambda}(\mathbb{P},\mathbb{Q}) = \sigma_{\mathbb{S},k,\lambda} \operatorname{MMD}_k(\mathbb{P},\mathbb{Q})$

Want to ensure
$$\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] \leq 1$$

Can solve with $\langle \partial_i \phi(x), f
angle_{\mathcal{H}} = \partial_i f(x)$...but too expensive!

Guaranteed if
$$\|f\|_{\mathcal{H}} \leq \sigma_{\mathbb{S},k,\lambda}$$

 $\sigma_{\mathbb{S},k,\lambda} := \left(\lambda + \mathop{\mathbb{E}}_{\tilde{X}\sim\mathbb{S}}\left[k(\tilde{X},\tilde{X}) + [\nabla_1 \cdot \nabla_2 k](\tilde{X},\tilde{X})\right]\right)^{-\frac{1}{2}}$

Gives distance $\mathrm{SMMD}_{\mathbb{S},k,\lambda}(\mathbb{P},\mathbb{Q}) = \sigma_{\mathbb{S},k,\lambda} \operatorname{MMD}_k(\mathbb{P},\mathbb{Q})$

$$egin{aligned} \mathcal{D}^{\Psi}_{ ext{MMD}} & ext{has} \ \mathcal{F} &= igcup_{\psi \in \Psi} \left\{ f: \|f\|_{\mathcal{H}_{\psi}} \ \leq 1
ight\} \ \mathcal{D}^{\mathbb{S},\Psi,\lambda}_{ ext{SMMD}} & ext{has} \ \mathcal{F} &= igcup_{\psi \in \Psi} \left\{ f: \|f\|_{\mathcal{H}_{\psi}} \ \leq \sigma_{\mathbb{S},k,\lambda}
ight\} \end{aligned}$$



$\mathop{\mathbb{E}}_{ ilde{X}\sim\mathbb{S}}[f(ilde{X})^2] + \mathop{\mathbb{E}}_{ ilde{X}\sim\mathbb{S}}[\| abla f(ilde{X})\|^2] + \lambda \|f\|_{\mathcal{H}}^2 \leq 1$

 $\mathbb{E}_{ ilde{X}\sim\mathbb{S}}[f(ilde{X})^2] + \mathbb{E}_{ ilde{X}\sim\mathbb{S}}[\|
abla f(ilde{X})\|^2] + \lambda \|f\|_{\mathcal{H}}^2 \leq 1 \ \mathbb{E}_{ ilde{X}\sim\mathbb{S}}[f(ilde{X})^2] = \left\langle f, \mathbb{E}_{ ilde{X}\sim\mathbb{S}}[k(ilde{X},\cdot)\otimes k(ilde{X},\cdot)]f
ight
angle$

$$\mathbb{E}_{ ilde{X} \sim \mathbb{S}} [f(ilde{X})^2] + \mathbb{E}_{ ilde{X} \sim \mathbb{S}} [\|
abla f(ilde{X})\|^2] + \lambda \|f\|_{\mathcal{H}}^2 \leq 1 \ \mathbb{E}_{ ilde{X} \sim \mathbb{S}} [f(ilde{X})^2] = \left\langle f, \mathbb{E}_{ ilde{X} \sim \mathbb{S}} [k(ilde{X}, \cdot) \otimes k(ilde{X}, \cdot)]f
ight
angle \ \mathbb{E}_{ ilde{X} \sim \mathbb{S}} [\|
abla f(ilde{X})\|^2] = \left\langle f, \mathbb{E}_{ ilde{X} \sim \mathbb{S}} \left[\sum_{i=1}^d \partial_i k(ilde{X}, \cdot) \otimes \partial_i k(ilde{X}, \cdot)
ight] f
ight
angle$$

$$\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[f(ilde{X})^2] + \mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] + \lambda \|f\|_{\mathcal{H}}^2 \leq 1$$
 $\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[f(ilde{X})^2] = \left\langle f, \mathbb{E}_{ ilde{X} \sim \mathbb{S}}\left[k(ilde{X}, \cdot) \otimes k(ilde{X}, \cdot)
ight]f
ight
angle$
 $\mathbb{E}_{ ilde{X} \sim \mathbb{S}}[\|
abla f(ilde{X})\|^2] = \left\langle f, \mathbb{E}_{ ilde{X} \sim \mathbb{S}}\left[\sum_{i=1}^d \partial_i k(ilde{X}, \cdot) \otimes \partial_i k(ilde{X}, \cdot)
ight]f
ight
angle$
 $\langle f, D_\lambda f
angle \leq \|D_\lambda\| \|f\|_{\mathcal{H}}^2 \leq \sigma_{\mathbb{S},k,\lambda}^{-2}\|f\|_{\mathcal{H}}^2$













Theorem: $\mathcal{D}_{\mathrm{SMMD}}^{\mathbb{S},\Psi,\lambda}$ is continuous.

If \mathbb{S} has a density; k_{top} is Gaussian/linear/...; ϕ_{ψ} is fully-connected, Leaky-ReLU, non-increasing width; all weights in Ψ have bounded condition number; then $\mathcal{W}(\mathbb{Q}_n, \mathbb{P}) \to 0$ implies $\mathcal{D}_{SMMD}^{\mathbb{S}, \Psi, \lambda}(\mathbb{Q}_n, \mathbb{P}) \to 0$.

Results on 160×160 CelebA

SN-SMMD-GAN

WGAN-GP





KID: 0.006

KID: 0.022



- Human evaluation: good at precision, bad at recall
- Likelihood: hard for GANs, maybe not right thing anyway
- Two-sample tests: always reject!

- Human evaluation: good at precision, bad at recall
- Likelihood: hard for GANs, maybe not right thing anyway
- Two-sample tests: always reject!
- Most common: Fréchet Inception Distance, FID
 - Run pretrained featurizer on model and target
 - Model each as Gaussian; compute W_2

- Human evaluation: good at precision, bad at recall
- Likelihood: hard for GANs, maybe not right thing anyway
- Two-sample tests: always reject!
- Most common: Fréchet Inception Distance, FID
 - Run pretrained featurizer on model and target
 - Model each as Gaussian; compute W_2
 - Strong bias, small variance: very misleading

- Human evaluation: good at precision, bad at recall
- Likelihood: hard for GANs, maybe not right thing anyway
- Two-sample tests: always reject!
- Most common: Fréchet Inception Distance, FID
 - Run pretrained featurizer on model and target
 - Model each as Gaussian; compute W_2
 - Strong bias, small variance: very misleading
 - Simple examples where $\operatorname{FID}(\mathbb{Q}_1) > \operatorname{FID}(\mathbb{Q}_2)$ but $\widehat{\operatorname{FID}}(\hat{\mathbb{Q}}_1) < \widehat{\operatorname{FID}}(\hat{\mathbb{Q}}_2)$ for reasonable sample size

- Human evaluation: good at precision, bad at recall
- Likelihood: hard for GANs, maybe not right thing anyway
- Two-sample tests: always reject!
- Most common: Fréchet Inception Distance, FID
 - Run pretrained featurizer on model and target
 - Model each as Gaussian; compute W_2
 - Strong bias, small variance: very misleading
 - Simple examples where $\operatorname{FID}(\mathbb{Q}_1) > \operatorname{FID}(\mathbb{Q}_2)$ but $\widehat{\operatorname{FID}}(\hat{\mathbb{Q}}_1) < \widehat{\operatorname{FID}}(\hat{\mathbb{Q}}_2)$ for reasonable sample size
- Our KID: MMD^2 instead. Unbiased, asymptotically normal



Training process on CelebA



Training process on CelebA $KID \times 10^{3}$ **SN-SMMDGAN** WGAN-GP MMDGAN-GP-L2 $\times 10^4$ generator iterations

Training process on CelebA



Training process on CelebA



IV: Unnormalized density/score estimation

- Problem: given samples $X_i \sim \mathbb{P}_0$ with density p_0
- Model is kernel exponential family: for any $f\in\mathcal{H}$,



IV: Unnormalized density/score estimation

- Problem: given samples $X_i \sim \mathbb{P}_0$ with density p_0
- Model is kernel exponential family: for any $f\in\mathcal{H}$,



i.e. any density with $\log p - \log q \in \mathcal{H}$

IV: Unnormalized density/score estimation

- Problem: given samples $X_i \sim \mathbb{P}_0$ with density p_0
- Model is kernel exponential family: for any $f\in\mathcal{H}$,



i.e. any density with $\log p - \log q \in \mathcal{H}$

• Gaussian k: dense in all continuous distributions on compact domains

Density estimation with KEFs

- Fitting with maximum likelihood is tough:
 - Z(f), abla Z(f) are tough to compute
 - Likelihood equations ill-posed for *characteristic* kernels
- We choose to fit the unnormalized model
 - Could then estimate Z(f) once after fitting if necessary

Unnormalized density / score estimation

- Don't necessarily need to compute Z(f) afterwards
- $f + \log q = \log p_f + \log Z(f)$, the "energy", lets us:
 - Find modes (global or local)
 - Sample (with MCMC)
 - ••••

. . .

- The score, $abla_x [f(x) + \log q(x)] =
 abla_x \log p_f(x)$, lets us:
 - Run HMC for targets whose gradients we can't evaluate
 - Construct Monte Carlo control functionals

Score matching in KEFs [Sriperumbudur+ JMLR-17]

• Idea: minimize Fisher divergence $J(p_0 \| p_f)$

$$J(f) = rac{1}{2} \int p_0(x) \|
abla_x \log p_f(x) -
abla_x \log p_0(x) \|^2 \mathrm{d}x$$

Score matching in KEFs [Sriperumbudur+ JMLR-17]

• Idea: minimize Fisher divergence $J(p_0 \| p_f)$

$$J(f) = rac{1}{2} \int p_0(x) \|
abla_x \log p_f(x) -
abla_x \log p_0(x) \|^2 \mathrm{d}x$$

• Under mild assumptions, $J(f) = C(p_0) + C(p_0)$

$$\int p_0(x) \sum_{d=1}^D \left[\partial_d^2 \log p_f(x) + rac{1}{2} (\partial_d \log p_f(x)))^2
ight] \mathrm{d}x$$

Score matching in KEFs [Sriperumbudur+ JMLR-17]

• Idea: minimize Fisher divergence $J(p_0 \| p_f)$

$$J(f) = rac{1}{2} \int p_0(x) \|
abla_x \log p_f(x) -
abla_x \log p_0(x) \|^2 \mathrm{d}x$$

• Under mild assumptions, $J(f) = C(p_0) + C(p_0)$

$$\int p_0(x) \sum_{d=1}^D \left[\partial_d^2 \log p_f(x) + rac{1}{2} (\partial_d \log p_f(x)))^2
ight] \mathrm{d}x$$

• Can estimate with Monte Carlo
• Minimize regularized loss function:

$$\hat{{J}}_{\lambda}(f) = rac{1}{n}\sum_{a=1}^{n}\sum_{i=1}^{d}\left[\partial_{i}^{2}f(X_{a}) + rac{1}{2}(\partial_{i}f(X_{a}))^{2}
ight] + rac{1}{2}\lambda\|f\|_{\mathcal{H}}^{2}$$

• Minimize regularized loss function:

$$\hat{{J}}_{\lambda}(f) = rac{1}{n}\sum_{a=1}^{n}\sum_{i=1}^{d}\left[\partial_{i}^{2}f(X_{a}) + rac{1}{2}(\partial_{i}f(X_{a}))^{2}
ight] + rac{1}{2}\lambda\|f\|_{\mathcal{H}}^{2}$$

- Representer theorem tells us minimizer of \hat{J}_{λ} over ${\cal H}$ is

$$f_{\lambda,\mathcal{X}}\in \mathrm{span}\left\{\partial_i k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}\cup \mathrm{span}\left\{\partial_i^2 k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}$$

• Minimize regularized loss function:

$$\hat{{J}}_{\lambda}(f) = rac{1}{n}\sum_{a=1}^{n}\sum_{i=1}^{d}\left[\partial_{i}^{2}f(X_{a}) + rac{1}{2}(\partial_{i}f(X_{a}))^{2}
ight] + rac{1}{2}\lambda\|f\|_{\mathcal{H}}^{2}$$

- Representer theorem tells us minimizer of \hat{J}_{λ} over ${\cal H}$ is

$$\{f_{\lambda,\mathcal{X}}\in \mathrm{span}\left\{\partial_i k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}\cup \mathrm{span}\left\{\partial_i^2 k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}$$



• Minimize regularized loss function:

$$\hat{{J}}_{\lambda}(f) = rac{1}{n}\sum_{a=1}^{n}\sum_{i=1}^{d}\left[\partial_{i}^{2}f(X_{a}) + rac{1}{2}(\partial_{i}f(X_{a}))^{2}
ight] + rac{1}{2}\lambda\|f\|_{\mathcal{H}}^{2}$$

- Representer theorem tells us minimizer of \hat{J}_{λ} over ${\cal H}$ is

$$f_{\lambda,\mathcal{X}}\in \mathrm{span}\left\{\partial_i k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}\cup \mathrm{span}\left\{\partial_i^2 k_{X_a}
ight\}_{a\in[n]}^{i\in[d]}$$

- Best $f\in \mathcal{H}$ is in

 $\mathcal{H}_{ ext{full}} = ext{span} \left\{ \partial_i k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]} \cup ext{span} \left\{ \partial_i^2 k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]}$

- Best $f\in \mathcal{H}$ is in

$$\mathcal{H}_{ ext{full}} = ext{span} \left\{ \partial_i k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]} \cup ext{span} \left\{ \partial_i^2 k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]}$$

• Find best f in dim M subspace in $\mathcal{O}(ndM^2+M^3)$ time

• Best $f\in \mathcal{H}$ is in

$$\mathcal{H}_{ ext{full}} = ext{span} \left\{ \partial_i k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]} \cup ext{span} \left\{ \partial_i^2 k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]}$$

• Find best f in dim M subspace in $\mathcal{O}(ndM^2+M^3)$ time

$$\beta = -(\frac{1}{n} \underbrace{B_{XY}^{\mathsf{T}}}_{M \times nd} \underbrace{B_{XY}}_{nd \times M} + \lambda \underbrace{G_{YY}}_{M \times M})^{\dagger} \underbrace{h_{Y}}_{M \times 1}$$

• Best $f\in \mathcal{H}$ is in

$$\mathcal{H}_{ ext{full}} = ext{span} \left\{ \partial_i k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]} \cup ext{span} \left\{ \partial_i^2 k_{X_a}
ight\}_{a \in [m{n}]}^{i \in [d]}$$

• Find best f in dim M subspace in $\mathcal{O}(ndM^2+M^3)$ time

$$\beta = -(\frac{1}{n} \underbrace{B_{XY}^{\mathsf{T}}}_{M \times nd} \underbrace{B_{XY}}_{nd \times M} + \lambda \underbrace{G_{YY}}_{M \times M})^{\dagger} \underbrace{h_{Y}}_{M \times 1}$$

• $\mathcal{H}_{\mathrm{full}}$: M=2nd, $\mathcal{O}(n^3d^3)$ time!

Nyström approximation [Sutherland+ AISTATS-18] $i \in [d]$

 $\bullet \; \mathcal{H}_{\mathrm{full}} = \mathrm{span} \left\{ \partial_i k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]} \cup \mathrm{span} \left\{ \partial_i^2 k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]}$



Nyström approximation [Sutherland+ AISTATS-18]

- $\bullet \; \mathcal{H}_{\mathrm{full}} = \mathrm{span} \left\{ \partial_i k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]} \cup \mathrm{span} \left\{ \partial_i^2 k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]}$
- Nyström approximation: find fit in different (smaller) \mathcal{H}_Y



Nyström approximation [Sutherland+ AISTATS-18]

- $\bullet \; \mathcal{H}_{\mathrm{full}} = \mathrm{span} \left\{ \partial_i k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]} \cup \mathrm{span} \left\{ \partial_i^2 k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]}$
- Nyström approximation: find fit in different (smaller) \mathcal{H}_Y
- One choice: pick $Y \subset [n]$, |Y| = m at random, then $\mathcal{H}_{nys}^Y = \operatorname{span} \{\partial_i k_{X_a}\}_{a \in Y}^{i \in [d]} \qquad \mathcal{O}(nm^2d^3) \text{ time}$ Get the same rates with $m = \sqrt{n} \log n$ (sometimes less)



Nyström approximation [Sutherland+ AISTATS-18]

- $\bullet \; \mathcal{H}_{\mathrm{full}} = \mathrm{span} \left\{ \partial_i k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]} \cup \mathrm{span} \left\{ \partial_i^2 k_{X_a} \right\}_{a \in [\boldsymbol{n}]}^{i \in [d]}$
- Nyström approximation: find fit in different (smaller) \mathcal{H}_Y
- One choice: pick $Y \subset [n]$, |Y| = m at random, then $\mathcal{H}_{nys}^Y = \operatorname{span} \{\partial_i k_{X_a}\}_{a \in Y}^{i \in [d]} \qquad \mathcal{O}(nm^2d^3) \text{ time}$ Get the same rates with $m = \sqrt{n} \log n$ (sometimes less)
- *"lite"*: pick Y at random, then $\mathcal{H}^Y_{ ext{lite}} = ext{span} \left\{ k_{X_a}
 ight\}_{a \in Y} \quad \mathcal{O}(nm^2d) ext{ time }$





Meta-learning a kernel

Meta-learning a kernel

Meta-learning a kernel

Results

• Learns local dataset geometry: better fits



• On real data: slightly worse likelihoods, maybe better "shapes" than deep likelihood models

Results

• Learns local dataset geometry: better fits



• On real data: slightly worse likelihoods, maybe better "shapes" than deep likelihood models

Results

• Learns local dataset geometry: better fits



• On real data: slightly worse likelihoods, maybe better "shapes" than deep likelihood models



Recap

Combining a deep architecture with a kernel machine that takes the higher-level learned representation as input can be quite powerful.

— Y. Bengio & Y. LeCun (2007), "Scaling Learning Algorithms towards Al"

Recap

Combining a deep architecture with a kernel machine that takes the higher-level learned representation as input can be quite powerful. — Y. Bengio & Y. LeCun (2007), "Scaling Learning Algorithms towards AI"

- Two-sample testing [ICLR-17, ICML-20, NeurIPS-21]
 - ψ maximizing power criterion, for one task or many
- Self-supervised learning with HSIC [NeurIPS-21]
 - Much better understanding of what's going on!
- Generative modeling with MMD GANs [ICLR-18, NeurIPS-18]
 Need a smooth loss function for the generator
- Score matching in exponential families [AISTATS-18, ICML-19]
 - Avoid overfitting with closed-form fit on held-out data