# Kernel distances between distributions for generative models

#### **Dougal J. Sutherland**

Gatsby Computational Neuroscience Unit, University College London

Michael Arbel Mikołaj Bińkowski Arthur Gretton UCL Imperial UCL



ICIAM 2019

Distance Metrics and Mass Transfer Between High Dimensional Point Clouds

 $\mathcal{D}_{\mathcal{F}}(\mathbb{P},\mathbb{Q}) = \sup_{f\in\mathcal{F}} \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[f(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[f(Y)]$ 

 $f:\mathcal{X}
ightarrow\mathbb{R}$  is a critic function

 $\mathcal{D}_\mathcal{F}(\mathbb{P},\mathbb{Q}) = \sup_{f\in\mathcal{F}} \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[f(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[f(Y)]$ 

 $f:\mathcal{X}
ightarrow\mathbb{R}$  is a critic function

$$\mathcal{D}_\mathcal{F}(\mathbb{P},\mathbb{Q}) = \sup_{f\in\mathcal{F}} \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[f(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[f(Y)]$$

 $f:\mathcal{X}
ightarrow\mathbb{R}$  is a critic function

Total variation:  $\mathcal{F} = \{f : f \text{ continuous}, |f(x)| \leq 1\}$ 



$$\mathcal{D}_{\mathcal{F}}(\mathbb{P},\mathbb{Q}) = \sup_{f\in\mathcal{F}} \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[f(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[f(Y)]$$

 $f:\mathcal{X}
ightarrow\mathbb{R}$  is a critic function

Total variation:  $\mathcal{F} = \{f: f \text{ continuous}, |f(x)| \leq 1\}$ Wasserstein:  $\mathcal{F} = \{f: \|f\|_{ ext{Lip}} \leq 1\}$ 



$$\mathcal{D}_{\mathcal{F}}(\mathbb{P},\mathbb{Q}) = \sup_{f\in\mathcal{F}} \mathop{\mathbb{E}}_{X\sim\mathbb{P}}[f(X)] - \mathop{\mathbb{E}}_{Y\sim\mathbb{Q}}[f(Y)]$$

 $f:\mathcal{X}
ightarrow\mathbb{R}$  is a critic function

Total variation:  $\mathcal{F} = \{f: f \text{ continuous}, |f(x)| \leq 1\}$ Wasserstein:  $\mathcal{F} = \{f: \|f\|_{ ext{Lip}} \leq 1\}$ 



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

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



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

$$f^*(t) \propto \mathop{\mathbb{E}}\limits_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}\limits_{Y \sim \mathbb{Q}} k(t,Y)$$



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

Kernel  $k:\mathcal{X} imes\mathcal{X} o\mathbb{R}$  – a "similarity" function

$$f^*(t) \propto \mathop{\mathbb{E}}_{X \sim \mathbb{P}} k(t,X) - \mathop{\mathbb{E}}_{Y \sim \mathbb{Q}} k(t,Y)$$

For many kernels,  $\mathrm{MMD}(\mathbb{P},\mathbb{Q})=0$  iff  $\mathbb{P}=\mathbb{Q}$ 



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

Reproducing property: if  $f\in \mathcal{H}_k$  ,  $f(x)=\langle f,arphi(x)
angle_{\mathcal{H}_k}$ 

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

Reproducing property: if  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}_k}$   $\operatorname{MMD}_k(\mathbb{P}, \mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[\langle f, \varphi(X) \rangle_{\mathcal{H}_k}] - \mathbb{E}_{Y \sim \mathbb{Q}}[\langle f, \varphi(Y) \rangle_{\mathcal{H}_k}]$ 

Reproducing property: if  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}_k}$  $\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[\langle f, arphi(X) 
angle_{\mathcal{H}_k}] - \mathbb{E}_{Y \sim \mathbb{Q}}[\langle f, arphi(Y) 
angle_{\mathcal{H}_k}]$  $= \sup_{\|f\|_{\mathcal{H}_L} \leq 1} \left\langle f, \mathop{\mathbb{E}}_{X \sim \mathbb{P}} [arphi(X)] - \mathop{\mathbb{E}}_{Y \sim \mathbb{Q}} [arphi(Y)] 
ight
angle_{\mathcal{H}_L} 
ight
angle$ 

Reproducing property: if  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}_k}$  $\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[\langle f, arphi(X) 
angle_{\mathcal{H}_k}] - \mathbb{E}_{Y \sim \mathbb{Q}}[\langle f, arphi(Y) 
angle_{\mathcal{H}_k}]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle f, \mathop{\mathbb{E}}_{X \sim \mathbb{P}} [arphi(X)] - \mathop{\mathbb{E}}_{Y \sim \mathbb{Q}} [arphi(Y)] 
ight
angle_{\mathcal{H}_k}$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle f, \mu^k_\mathbb{P} - \mu^k_\mathbb{Q} 
ight
angle_{\mathcal{H}_k}$ 

Reproducing property: if  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}_k}$  $\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[\langle f, arphi(X) 
angle_{\mathcal{H}_k}] - \mathbb{E}_{Y \sim \mathbb{Q}}[\langle f, arphi(Y) 
angle_{\mathcal{H}_k}]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle f, \mathop{\mathbb{E}}_{X \sim \mathbb{P}} [arphi(X)] - \mathop{\mathbb{E}}_{Y \sim \mathbb{Q}} [arphi(Y)] 
ight
angle_{\mathcal{H}_k}$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle f, \mu^k_{\mathbb{P}} - \mu^k_{\mathbb{Q}} 
ight
angle_{\mathcal{H}_k} = \left\| \mu^k_{\mathbb{P}} - \mu^k_{\mathbb{Q}} 
ight\|_{\mathcal{H}_k}$ 

Reproducing property: if  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}_k}$  $\mathrm{MMD}_k(\mathbb{P},\mathbb{Q}) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[\langle f, arphi(X) 
angle_{\mathcal{H}_k}] - \mathbb{E}_{Y \sim \mathbb{Q}}[\langle f, arphi(Y) 
angle_{\mathcal{H}_k}]$  $= \sup_{\|f\|_{\mathcal{H}_L} \leq 1} \left\langle f, \mathop{\mathbb{E}}_{X \sim \mathbb{P}} [arphi(X)] - \mathop{\mathbb{E}}_{Y \sim \mathbb{Q}} [arphi(Y)] 
ight
angle_{\mathcal{H}_L}$  $= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle f, \mu^k_{\mathbb{P}} - \mu^k_{\mathbb{Q}} 
ight
angle_{\mathcal{H}_k} = \left\| \mu^k_{\mathbb{P}} - \mu^k_{\mathbb{Q}} 
ight\|_{\mathcal{H}_k}$  $\langle \mu^k_{\mathbb{P}}, \mu^k_{\mathbb{Q}} 
angle_{\mathcal{H}_k} = \mathop{\mathbb{E}}_{\substack{X \sim \mathbb{P} \ Y \sim \mathbb{Q}}} \langle \varphi(X), \varphi(Y) 
angle_{\mathcal{H}_k} = \mathop{\mathbb{E}}_{\substack{X \sim \mathbb{P} \ Y \sim \mathbb{Q}}} k(X, Y)$ 

#### MMD as feature matching

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

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

#### MMD as feature matching

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

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

#### MMD as feature matching

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

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

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

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

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





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

 $K_{XX}$ 



 $K_{YY}$ 



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

 $K_{XX}$ 











# Implicit generative models

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



 $X \sim \mathbb{P}$ 



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

# Why implicit generative models?

# Why implicit generative models?



VS



 $X\sim \mathbb{P}$ 

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

# Why implicit generative models?

#### Uses of generative models:

• Automated animation (anime characters)



 $X \sim \mathbb{P}$ 



 $Y \sim \mathbb{Q}_A$


#### Uses of generative models:

• Automated animation (anime characters)



 $X \sim \mathbb{P}$ 



 $Y \sim \mathbb{Q}_A$ 

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)



 $X\sim \mathbb{P}$ 



 $Y \sim \mathbb{O}_{A}$ 

Automated	字	種	成	東	字	推			
• Learn callig	符	利	對	亞	型	斷			
	到	用	抗	話	進	的			
	字	條	網		行	新			
	符	件	絡	字	自	方			
		生	對	體	動	法			
	1	Ш			-	K B	<b>,</b>		

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)



 $X\sim \mathbb{P}$ 



 $Y \sim \mathbb{O}_{A}$ 

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)
- Image translation (pix2pix, Wolf+ 17, Everybody Dance Now)



 $X\sim \mathbb{P}$ 



 $Y \sim \mathbb{O}$ 

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)







#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)
- Image translation (pix2pix, Wolf+ 17, Everybody Dance Now)



 $X\sim \mathbb{P}$ 



 $Y \sim \mathbb{O}$ 

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)
- Image translation (pix2pix, Wolf+ 17, Everybody Dance Now)
- Help musicians improvise, "natural" image editing, plan robot actions, anonymize datasets, ...



 $X \sim \mathbb{P}$ 



#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)
- Image translation (pix2pix, Wolf+ 17, Everybody Dance Now)
- Help musicians improvise, "natural" image editing, plan robot actions, anonymize datasets, ...
- Make \$432,500 selling a generated image





 $X \sim \mathbb{P}$ 



#### Why implicit apparative models?

- Automated
- Learn callig
- Image trans
- Help musici plan robot
- Make \$432,





# Is artificial intelligence set to become art's next medium?

Al artwork sells for \$432,500 — nearly 45 times its high estimate — as Christie's becomes the first auction house to offer a work of art created by an algorithm

The portrait in its gilt frame depicts a portly gentleman, possibly French and — to judge by his dark frockcoat and plain white collar — a man of the church. The work appears unfinished: the facial features are somewhat indistinct and there are blank areas of canvas. Oddly, the whole composition is displaced slightly to the north-west. A label on the wall states that the sitter is a man named Edmond Belamy, but the giveaway clue as to the origins of the work is the artist's signature at the bottom right. In cursive Gallic script it reads:

 $\min_{G} \max_{D} \mathbb{E}_{x}[\log(D(x))] + \mathbb{E}_{z}[\log(1 - D(G(z)))]$ 

Image © Obvious

y Dance Now)

ling,

#### Uses of generative models:

- Automated animation (anime characters)
- Learn calligraphic font style (zi2zi)
- Image translation (pix2pix, Wolf+ 17, Everybody Dance Now)
- Help musicians improvise, "natural" image editing, plan robot actions, anonymize datasets, ...
- Make \$432,500 selling a generated image





 $X \sim \mathbb{P}$ 



#### Model: generator network

Fixed distribution of latents:  $Z \sim \text{Uniform}\left([-1,1]^{100}
ight)$ 

Maps through a network:  $G_{m{ heta}}(Z) \sim \mathbb{Q}_{m{ heta}}$ 



#### Model: generator network

Fixed distribution of latents:  $Z \sim \text{Uniform}([-1, 1]^{100})$ 

Maps through a network:  $G_{m{ heta}}(Z) \sim \mathbb{Q}_{m{ heta}}$ 



```
Choose 	heta to minimize \mathcal{D}(\mathbb{P}_{\mathrm{data}},\mathbb{Q}_{	heta})
```

## Model: generator network

Fixed distribution of latents:  $Z \sim \text{Uniform}([-1, 1]^{100})$ 

Maps through a network:  $G_{m{ heta}}(Z) \sim \mathbb{Q}_{m{ heta}}$ 



```
Choose \theta to minimize \mathcal{D}(\mathbb{P}_{data}, \mathbb{Q}_{\theta})
Very flexible generative model
```

# Traditional choices for ${\cal D}$

 $\operatorname{argmin}_{\theta} \operatorname{KL}(\mathbb{P}_{\operatorname{data}} \| \mathbb{Q}_{\theta})$ 

maximum likelihood

 $\operatorname{argmin}_{\theta} \operatorname{JS}(\mathbb{P}_{\operatorname{data}}, \mathbb{Q}_{\theta})$ 

original GANs (ish)

#### 

• Problem:  $\mathbb{P}_{data}$  and  $\mathbb{Q}_{\theta}$  don't have full-dimensional support



# $\begin{aligned} & \text{Traditional choices for } \mathcal{D} \\ & \arg\min_{\theta} \mathrm{KL}(\mathbb{P}_{\mathrm{data}} \| \mathbb{Q}_{\theta}) \qquad \arg\min_{\theta} \mathrm{JS}(\mathbb{P}_{\mathrm{data}}, \mathbb{Q}_{\theta}) \end{aligned}$

original GANs (ish)

- Problem:  $\mathbb{P}_{data}$  and  $\mathbb{Q}_{\theta}$  don't have full-dimensional support
- Supports almost surely don't align [Arjovsky/Bottou ICLR-17]



maximum likelihood

## Traditional choices for ${\cal D}$



- Problem:  $\mathbb{P}_{data}$  and  $\mathbb{Q}_{\theta}$  don't have full-dimensional support
- Supports almost surely don't align [Arjovsky/Bottou ICLR-17]
- No "partial credit": loss is flat at maximum, never improves



#### MMD as loss [Li+ ICML-15, Dziugaite+ UAI-15]

- Does give "partial credit" to nearby points
- Can directly minimize  $\widehat{\mathrm{MMD}}^2(X,G_{{m heta}}(Z))$  with SGD

#### MMD as loss [Li+ ICML-15, Dziugaite+ UAI-15]

- Does give "partial credit" to nearby points
- Can directly minimize  $\widehat{\mathrm{MMD}}^2(X,G_{ heta}(Z))$  with SGD

MNIST, mix of Gaussian kernels





#### MMD as loss [Li+ ICML-15, Dziugaite+ UAI-15]

- Does give "partial credit" to nearby points
- Can directly minimize  $\widehat{\mathrm{MMD}}^2(X,G_{ heta}(Z))$  with SGD

#### MNIST, mix of Gaussian kernels





#### Celeb-A, mix of rational quadratic + linear kernels



 $\mathbb{P}_{\text{data}}$ 

<sup>II</sup> data

ĽĤ

#### Celeb-A, mix of rational quadratic + linear kernels



#### **Deep kernels**

$$egin{aligned} k_\psi(x,y) &= k_{ ext{top}}(\phi_\psi(x),\phi_\psi(y)) \ \phi_\psi:\mathcal{X} & o \mathbb{R}^D \qquad k_\psi:\mathbb{R}^D imes \mathbb{R}^D o \mathbb{R} \end{aligned}$$



•  $k_{
m top}$  usually Gaussian, linear, ...

# MMD loss with a deep kernel

$$k(x,y) = k_{ ext{top}}(\phi(x),\phi(y))$$

- $\phi: \mathcal{X} 
  ightarrow \mathbb{R}^{2048}$  from pretrained Inception net
- $k_{
  m top}$  simple: exponentiated quadratic or polynomial

# MMD loss with a deep kernel

$$k(x,y) = k_{ ext{top}}(\phi(x),\phi(y))$$

- $\phi: \mathcal{X} 
  ightarrow \mathbb{R}^{2048}$  from pretrained Inception net
- $k_{
  m top}$  simple: exponentiated quadratic or polynomial



 $\mathbb{P}_{data}$ 

# MMD loss with a deep kernel

$$k(x,y) = k_{ ext{top}}(\phi(x),\phi(y))$$

- $\phi: \mathcal{X} 
  ightarrow \mathbb{R}^{2048}$  from pretrained Inception net
- $k_{
  m top}$  simple: exponentiated quadratic or polynomial



nită.	Æ	đ.	
	 	Æ.	



# **Optimized MMD: MMD GANs [Li+ NeurIPS-17]**

• Don't just use one kernel, use a *class* parameterized by  $\psi$ :

$$k_\psi(x,y) = k_{ ext{top}}(\phi_\psi(x),\phi_\psi(y))$$

# **Optimized MMD: MMD GANs [Li+ NeurIPS-17]**

• Don't just use one kernel, use a *class* parameterized by  $\psi$ :

$$k_\psi(x,y) = k_{ ext{top}}(\phi_\psi(x),\phi_\psi(y))$$

• New distance based on *all* these kernels:

$$\mathcal{D}^{\Psi}_{\mathrm{MMD}}(\mathbb{P},\mathbb{Q}) = \sup_{\psi\in\Psi}\mathrm{MMD}_{\psi}(\mathbb{P},\mathbb{Q})$$

# **Optimized MMD: MMD GANs [Li+ NeurIPS-17]**

• Don't just use one kernel, use a *class* parameterized by  $\psi$ :

$$k_\psi(x,y) = k_{ ext{top}}(\phi_\psi(x),\phi_\psi(y))$$

• New distance based on *all* these kernels:

$$\mathcal{D}^{\Psi}_{\mathrm{MMD}}(\mathbb{P},\mathbb{Q}) = \sup_{\psi\in\Psi}\mathrm{MMD}_{\psi}(\mathbb{P},\mathbb{Q})$$

• Minimax optimization problem

$$\inf_{\substack{ heta \ \psi}} \operatorname{MMD}_{\psi}(\mathbb{P}_{\operatorname{data}}, \mathbb{Q}_{ heta})$$

# Non-smoothness of Optimized MMD

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



# Non-smoothness of Optimized MMD

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






























- ...deep kernel analogue is hard.
- Instead, keep witness function from being too steep

$$-5.0 -2.5 0.0 2.5 5.0 0 5 10 \theta$$



- Just need to stay away from tiny bandwidths  $\psi$
- ...deep kernel analogue is hard.
- Instead, keep witness function from being too steep
- Control  $\|
  abla f( ilde X)\|$  on average, near the data
  - [Gulrajani+ NeurIPS-17 / Roth+ NeurIPS-17 / Mescheder+ ICML-18]

-5.0	-2.5	0.0	2.5	5.0	0	5	10
		θ				$\theta$	

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

- If  $\Psi$  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  $\Psi$  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  $\Psi$  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 do directly with kernel properties...but too expensive!

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

Can do directly with kernel properties...but too expensive!

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

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

Can do directly with kernel properties...but too expensive!

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

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

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

Can do directly with kernel properties...but too expensive!

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

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

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

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

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

$$egin{aligned} &\mathbb{E}_{ ilde{X}\sim\mathbb{S}}[f( ilde{X})^2] + \mathbb{E}_{ ilde{X}\sim\mathbb{S}}[\|
abla f( ilde{X})\|^2] + \lambda \|f\|_{\mathcal{H}}^2 \leq 1 \ &\mathbb{E}_{ ilde{X}\sim\mathbb{S}}[f( ilde{X})^2] = \left\langle f, \mathbb{E}_{ ilde{X}\sim\mathbb{S}}[arphi( ilde{X})\otimesarphi( ilde{X})]f
ight
angle_{\mathcal{H}} \end{aligned}$$

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

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

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

$$\langle f, D_\lambda f 
angle_{\mathcal{H}}$$

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

$$\langle f, D_\lambda f 
angle_{\mathcal{H}} \leq \| D_\lambda \| \, \| f \|_{\mathcal{H}}^2$$

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

$$\langle f, D_\lambda f 
angle_{\mathcal{H}} \leq \| D_\lambda \| \, \| f \|_{\mathcal{H}}^2 \leq \sigma_{\mathbb{S},k,\lambda}^{-2} \| f \|_{\mathcal{H}}^2$$








# Smoothness of $\mathcal{D}_{\mathrm{SMMD}}$



## Smoothness of $\mathcal{D}_{\mathrm{SMMD}}$



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

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

#### Target $\mathbb{P}$ and model $\mathbb{Q}_{\theta}$ samples



Kernels from  $\mathrm{SMMD}_{\mathbb{P},k,\lambda}$ , early in optimization





## Critic gradients from $\mathrm{SMMD}_{\mathbb{P},k,\lambda}$ (early)

SMMDGAN (target)



#### Critic gradients from $\mathrm{MMD}_k$ (early)

MMDGAN (no GP)



#### Kernels from $\mathrm{SMMD}_{\mathbb{P},k,\lambda}$ , late in optimization

SMMDGAN (target)





### Critic gradients from $\mathrm{SMMD}_{\mathbb{P},k,\lambda}$ (late)



#### Critic gradients from $\mathrm{MMD}_k$ (late)



## Model on $160\times 160$ CelebA

#### SN-SMMD-GAN

#### WGAN-GP





#### KID: 0.006

#### KID: 0.022











## $\textbf{Model on } 64 \times 64 \textbf{ ImageNet}$

#### SN-SMMDGAN



**SN-GAN** 

BGAN





KID: 0.047

KID: 0.035

#### KID: 0.044

## Recap

- Can train generative models by minimizing a *flexible*, *smooth* distance between distributions
- Combine kernels with gradient penalties
- Strong practical results, some understanding of theory

#### Demystifying MMD GANs

Bińkowski<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Arbel, and Gretton [ICLR 2018]

#### **On Gradient Regularizers for MMD GANs**

Arbel<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Bińkowski, and Gretton [NeurIPS 2018]

Links + code: see dougal.me

## Recap

- Can train generative models by minimizing a *flexible*, *smooth* distance between distributions
- Combine kernels with gradient penalties
- Strong practical results, some understanding of theory
- But haven't totally solved GANs yet!

#### **Demystifying MMD GANs**

Bińkowski<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Arbel, and Gretton [ICLR 2018]

#### **On Gradient Regularizers for MMD GANs**

Arbel<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Bińkowski, and Gretton [NeurIPS 2018]

Links + code: see dougal.me

## Recap

- Can train generative models by minimizing a *flexible*, *smooth* distance between distributions
- Combine kernels with gradient penalties
- Strong practical results, some understanding of theory
- But haven't totally solved GANs yet!

#### **Demystifying MMD GANs**

Bińkowski<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Arbel, and Gretton [ICLR 2018]

#### **On Gradient Regularizers for MMD GANs**

Arbel<sup>\*</sup>, <u>Sutherland</u><sup>\*</sup>, Bińkowski, and Gretton [NeurIPS 2018]

Links + code: see dougal.me

Thanks!

## **Backup slides**

## **MMD GANs versus WGANs**

- Linear-kernel MMD GAN,  $k(x, y) = \phi(x)\phi(y)$ :  $\log x = |\underset{\mathbb{P}}{\mathbb{E}}\phi(X) - \underset{\mathbb{Q}}{\mathbb{E}}\phi(Y)|$  $f(t) = \operatorname{sign}\left(\underset{\mathbb{P}}{\mathbb{E}}\phi(X) - \underset{\mathbb{Q}}{\mathbb{E}}\phi(Y)\right)\phi(t)$
- WGAN has:

$$egin{aligned} &\mathrm{loss} = \mathop{\mathbb{E}}\limits_{\mathbb{P}} \phi(X) - \mathop{\mathbb{E}}\limits_{\mathbb{Q}} \phi(Y) \ &f(t) = \phi(t) \end{aligned}$$

## **MMD GANs versus WGANs**

- Linear-kernel MMD GAN,  $k(x, y) = \phi(x)\phi(y)$ :  $\log x = |\underset{\mathbb{P}}{\mathbb{E}} \phi(X) - \underset{\mathbb{Q}}{\mathbb{E}} \phi(Y)|$  $f(t) = \operatorname{sign}\left(\underset{\mathbb{P}}{\mathbb{E}} \phi(X) - \underset{\mathbb{Q}}{\mathbb{E}} \phi(Y)\right)\phi(t)$
- WGAN has:

$$egin{aligned} & \log & = \mathop{\mathbb{E}}\limits_{\mathbb{P}} \phi(X) - \mathop{\mathbb{E}}\limits_{\mathbb{Q}} \phi(Y) \ & f(t) = \phi(t) \end{aligned}$$

- Linear-kernel MMD GAN-GP and WGAN-GP almost the same
- MMD GAN "offloads" some of the critic's work to closed-form optimization in the RKHS

## **Keeping weight condition numbers bounded**

- Spectral parameterization [Miyato+ ICLR-18]:
- $W=\gamma ar{W}/\|ar{W}\|_{
  m op}$ ; learn  $\gamma$  and  $ar{W}$  freely
- Encourages diversity without limiting representation



## Rank collapse

- Occasional optimization failure without spectral param:
  - Generator doing reasonably well
  - Critic filters become low-rank
  - Generator corrects it by breaking everything else
  - Generator gets stuck



# What if we just did spectral normalization?

- $W=ar{W}/\|ar{W}\|_{ ext{op}}$  , so that  $\|W\|_{ ext{op}}=1$  ,  $\|\phi_{\psi}\|_{L}\leq 1$
- Works well for original GANs [Miyato+ ICLR-18]
- ...but doesn't work at all as only constraint in a WGAN
- Limits representation too much
  - In DiracGAN, only allows bandwidth 1
  - $\|x\mapsto \sigma(W_n\cdots\sigma(W_1x))\|_L \ll \|W_n\|_{\mathrm{op}}\cdots\|W_1\|_{\mathrm{op}}$

# • $k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$ means $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  
 $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

• Can show  $\mathrm{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\mathrm{top}}} \| \phi_\psi \|_{\mathrm{Lip}} \, \mathcal{W}$ 

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  
 $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

- Can show  $\mathrm{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\mathrm{top}}} \| \phi_\psi \|_{\mathrm{Lip}} \, \mathcal{W}$
- By assumption on  $k_{ ext{top}}$  ,  $\sigma_{\mathbb{S},k,\lambda}^{-2} \geq \gamma_{k_{ ext{top}}}^2 \ \mathbb{E}[\|
  abla \phi_\psi( ilde X)\|_F^2]$

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  
 $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

- Can show  $\operatorname{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\operatorname{top}}} \| \phi_\psi \|_{\operatorname{Lip}} \, \mathcal{W}$
- By assumption on  $k_{ ext{top}}$  ,  $\sigma_{\mathbb{S},k,\lambda}^{-2} \geq \gamma_{k_{ ext{top}}}^2 \ \mathbb{E}[\|
  abla \phi_\psi( ilde X)\|_F^2]$

$$\bullet \; \mathrm{SMMD}^2 \leq \frac{L^2_{k_{\mathrm{top}}} \| \phi_\psi \|^2_{\mathrm{Lip}}}{\gamma^2_{k_{\mathrm{top}}} \, \mathbb{E} \, \| \nabla_{\tilde{X}} \phi_\psi(\tilde{X}) \|^2_F} \, \mathcal{W}$$

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  
 $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

- Can show  $\operatorname{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\operatorname{top}}} \| \phi_\psi \|_{\operatorname{Lip}} \, \mathcal{W}$
- By assumption on  $k_{ ext{top}}$  ,  $\sigma_{\mathbb{S},k,\lambda}^{-2} \geq \gamma_{k_{ ext{top}}}^2 \ \mathbb{E}[\|
  abla \phi_\psi( ilde X)\|_F^2]$

$$\bullet \; \mathrm{SMMD}^2 \leq \frac{L^2_{k_{\mathrm{top}}} \| \phi_\psi \|^2_{\mathrm{Lip}}}{\gamma^2_{k_{\mathrm{top}}} \; \mathbb{E} \, \| \nabla_{\tilde{X}} \phi_\psi(\tilde{X}) \|_F^2} \, \mathcal{W}$$

- Because Leaky-ReLU,  $\phi_\psi(X) = lpha(\psi) \phi_{ar\psi}(X)$ ,  $\|\phi_{ar\psi}\|_{ ext{Lip}} \leq 1$ 

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

- Can show  $\operatorname{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\operatorname{top}}} \| \phi_\psi \|_{\operatorname{Lip}} \, \mathcal{W}$
- By assumption on  $k_{ ext{top}}$  ,  $\sigma_{\mathbb{S},k,\lambda}^{-2} \geq \gamma_{k_{ ext{top}}}^2 \ \mathbb{E}[\|
  abla \phi_\psi( ilde X)\|_F^2]$

$$\bullet \; \mathrm{SMMD}^2 \leq \frac{L^2_{k_{\mathrm{top}}} \| \phi_\psi \|^2_{\mathrm{Lip}}}{\gamma^2_{k_{\mathrm{top}}} \, \mathbb{E} \, \| \nabla_{\tilde{X}} \phi_\psi(\tilde{X}) \|_F^2} \, \mathcal{W}$$

- Because Leaky-ReLU,  $\phi_\psi(X) = lpha(\psi) \phi_{ar\psi}(X)$ ,  $\|\phi_{ar\psi}\|_{ ext{Lip}} \leq 1$
- For Lebesgue-almost all  $ilde{X}$ ,  $\| 
  abla_{ ilde{X}} \phi_{ar{\psi}}( ilde{X}) \|_F^2 \geq rac{d_{ ext{top}} lpha^L}{\kappa^L}$

• 
$$k_{\psi}(x,y) = k_{ ext{top}}(\phi_{\psi}(x),\phi(y))$$
 means  $d_{\psi}(x,y) = \|k_{\psi}(x,\cdot) - k_{\psi}(y,\cdot)\|_{\mathcal{H}_{k_{\psi}}} \leq L_{k_{ ext{top}}}\|\phi_{\psi}\|_{ ext{Lip}}\|x-y\|$ 

- Can show  $\operatorname{MMD}_\psi \leq \mathcal{W}_{d_\psi} \leq L_{k_{\operatorname{top}}} \| \phi_\psi \|_{\operatorname{Lip}} \, \mathcal{W}$
- By assumption on  $k_{ ext{top}}$  ,  $\sigma_{\mathbb{S},k,\lambda}^{-2} \geq \gamma_{k_{ ext{top}}}^2 \ \mathbb{E}[\|
  abla \phi_\psi( ilde X)\|_F^2]$
- $\bullet \ \mathrm{SMMD}^2 \leq \frac{L_{k_{\mathrm{top}}}^2 \|\phi_\psi\|_{\mathrm{Lip}}^2}{\gamma_{k_{\mathrm{top}}}^2 \, \mathbb{E} \, \|\nabla_{\tilde{X}} \phi_\psi(\tilde{X})\|_F^2} \, \mathcal{W} \leq \frac{L_{k_{\mathrm{top}}}^2 \, \kappa^L}{\gamma_{k_{\mathrm{top}}}^2 \, d_{\mathrm{top}} \alpha^L} \, \mathcal{W}$
- Because Leaky-ReLU,  $\phi_\psi(X) = lpha(\psi) \phi_{ar\psi}(X)$ ,  $\|\phi_{ar\psi}\|_{ ext{Lip}} \leq 1$
- For Lebesgue-almost all  $ilde{X}$ ,  $\| 
  abla_{ ilde{X}} \phi_{ar{\psi}}( ilde{X}) \|_F^2 \geq rac{d_{ ext{top}} lpha^L}{\kappa^L}$

## Implicit generative model evaluation

• No likelihoods, so...how to compare models?

## Implicit generative model evaluation

- No likelihoods, so...how to compare models?
- Main approach:

look at a bunch of pictures and see if they're pretty or not
- No likelihoods, so...how to compare models?
- Main approach:
  - look at a bunch of pictures and see if they're pretty or not
    - Easy to find (really) bad samples

- No likelihoods, so...how to compare models?
- Main approach:
  - look at a bunch of pictures and see if they're pretty or not
    - Easy to find (really) bad samples
    - Hard to see if modes are missing / have wrong probabilities

- No likelihoods, so...how to compare models?
- Main approach:
  - look at a bunch of pictures and see if they're pretty or not
    - Easy to find (really) bad samples
    - Hard to see if modes are missing / have wrong probabilities
    - Hard to compare models beyond certain threshold

- No likelihoods, so...how to compare models?
- Main approach:
  - look at a bunch of pictures and see if they're pretty or not
    - Easy to find (really) bad samples
    - Hard to see if modes are missing / have wrong probabilities
    - Hard to compare models beyond certain threshold
- Need better, quantitative methods

- No likelihoods, so...how to compare models?
- Main approach: look at a bunch of pictures and see if they're pretty or not
  - Easy to find (really) bad samples
  - Hard to see if modes are missing / have wrong probabilities
  - Hard to compare models beyond certain threshold
- Need better, quantitative methods
- Our method: Kernel Inception Distance (KID)
  - $\mathrm{MMD}_k(\mathbb{P}_{\mathrm{data}},\mathbb{Q}_{ heta})^2$ , k cubic on pretrained Inception rep

- No likelihoods, so...how to compare models?
- Main approach: look at a bunch of pictures and see if they're pretty or not
  - Easy to find (really) bad samples
  - Hard to see if modes are missing / have wrong probabilities
  - Hard to compare models beyond certain threshold
- Need better, quantitative methods
- Our method: Kernel Inception Distance (KID)
  - $\mathrm{MMD}_k(\mathbb{P}_{\mathrm{data}},\mathbb{Q}_{ heta})^2$ , k cubic on pretrained Inception rep
  - tf.contrib.gan.eval.kernel\_inception\_distance

- Previously standard quantitative method
- Based on ImageNet classifier label predictions
  - Classifier should be confident on individual images
  - Predicted labels should be diverse across sample

- Previously standard quantitative method
- Based on ImageNet classifier label predictions
  - Classifier should be confident on individual images
  - Predicted labels should be diverse across sample
- No notion of target distribution  $\mathbb{P}_{data}$

- Previously standard quantitative method
- Based on ImageNet classifier label predictions
  - Classifier should be confident on individual images
  - Predicted labels should be diverse across sample
- No notion of target distribution  $\mathbb{P}_{data}$
- Scores completely meaningless on LSUN, Celeb-A, SVHN, ...

- Previously standard quantitative method
- Based on ImageNet classifier label predictions
  - Classifier should be confident on individual images
  - Predicted labels should be diverse across sample
- No notion of target distribution  $\mathbb{P}_{data}$
- Scores completely meaningless on LSUN, Celeb-A, SVHN, ...
- Not great on CIFAR-10 either

• Previously standard quantitative method



- Fit normals to Inception hidden layer activations of  ${\mathbb P}$  and  ${\mathbb Q}$
- Compute Fréchet (Wasserstein-2) distance between fits

- Fit normals to Inception hidden layer activations of  ${\mathbb P}$  and  ${\mathbb Q}$
- Compute Fréchet (Wasserstein-2) distance between fits
- Meaningful on not-ImageNet datasets

- Fit normals to Inception hidden layer activations of  $\mathbb P$  and  $\mathbb Q$
- Compute Fréchet (Wasserstein-2) distance between fits
- Meaningful on not-ImageNet datasets



- Fit normals to Inception hidden layer activations of  ${\mathbb P}$  and  ${\mathbb Q}$
- Compute Fréchet (Wasserstein-2) distance between fits
- Meaningful on not-ImageNet datasets
- Estimator extremely biased, tiny variance



- Fit normals to Inception hidden layer activations of  ${\mathbb P}$  and  ${\mathbb Q}$
- Compute Fréchet (Wasserstein-2) distance between fits
- Meaningful on not-ImageNet datasets
- Estimator extremely biased, tiny variance
- $\operatorname{FID}(\mathbb{P}_1,\mathbb{Q}) < \operatorname{FID}(\mathbb{P}_2,\mathbb{Q}), \mathbb{E}\operatorname{FID}(\hat{\mathbb{P}}_1,\mathbb{Q}) > \mathbb{E}\operatorname{FID}(\hat{\mathbb{P}}_2,\mathbb{Q})$





#### New method: Kernel Inception Distance (KID)

- $\widehat{\mathrm{MMD}}^2$  between Inception hidden layer activations
- Use default polynomial kernel:  $k(x,y) = \left(rac{1}{d} \langle x,y 
  angle + 1
  ight)^3$

#### New method: Kernel Inception Distance (KID)

- $\widehat{\mathrm{MMD}}^2$  between Inception hidden layer activations
- Use default polynomial kernel:  $k(x,y) = \left(rac{1}{d} \langle x,y 
  angle + 1
  ight)^3$
- Unbiased estimator, reasonable with few samples



#### New method: Kernel Inception Distance (KID)

- $\widehat{\mathrm{MMD}}^2$  between Inception hidden layer activations
- Use default polynomial kernel:  $k(x,y) = \left(rac{1}{d} \langle x,y 
  angle + 1
  ight)^3$
- Unbiased estimator, reasonable with few samples
- In tensorflow.contrib.gan.eval (tensorflow#21066)



### Automatic learning rate adaptation with KID

- Models need appropriate learning rate schedule to work well
- Automate with three-sample MMD test [Bounliphone+ ICLR-16]:



### **Controlling critic complexity**

