

Generative Adversarial Nets and Distances Between Distributions



(from thispersondoesnotexist.com)

CPSC 340, Fall 2021

Generative Adversarial Nets and Distances Between Distributions



(from thispersondoesnotexist.com)

CPSC 340, Fall 2021

(we definitely won't get through all these slides, but let's see!)

Generative models

- Start with a bunch of examples: $\mathcal{X}_1, \dots, \mathcal{X}_n \sim \mathbb{P}$
- Want a model for the data: $\mathcal{Q} \approx \mathbb{P}$

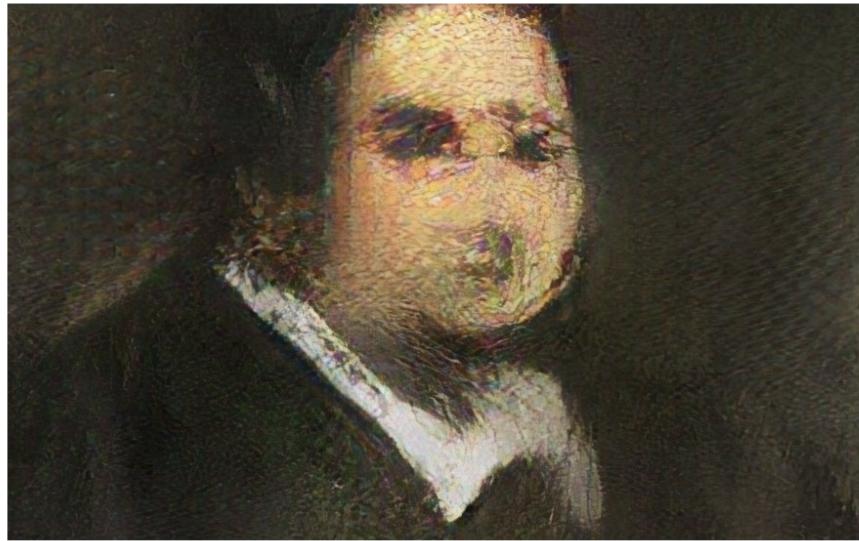
Generative models

- Start with a bunch of examples: $\mathbf{X}_1, \dots, \mathbf{X}_n \sim \mathbb{P}$
- Want a model for the data: $\mathbb{Q} \approx \mathbb{P}$
- Might want to do different things with the model:
 - Find most representative data points / modes
 - Find outliers, anomalies, ...
 - Discover underlying structure of the data
 - Impute missing values
 - Use as prior (semi-supervised, machine translation, ...)
 - Produce “more samples”
 - ...

Generative models

- Start with a bunch of examples: $\mathbf{X}_1, \dots, \mathbf{X}_n \sim \mathbb{P}$
- Want a model for the data: $\mathbb{Q} \approx \mathbb{P}$
- Might want to do different things with the model:
 - Find most representative data points / modes
 - Find outliers, anomalies, ...
 - Discover underlying structure of the data
 - Impute missing values
 - Use as prior (semi-supervised, machine translation, ...)
 - Produce “more samples”
 - ...

Why produce samples?



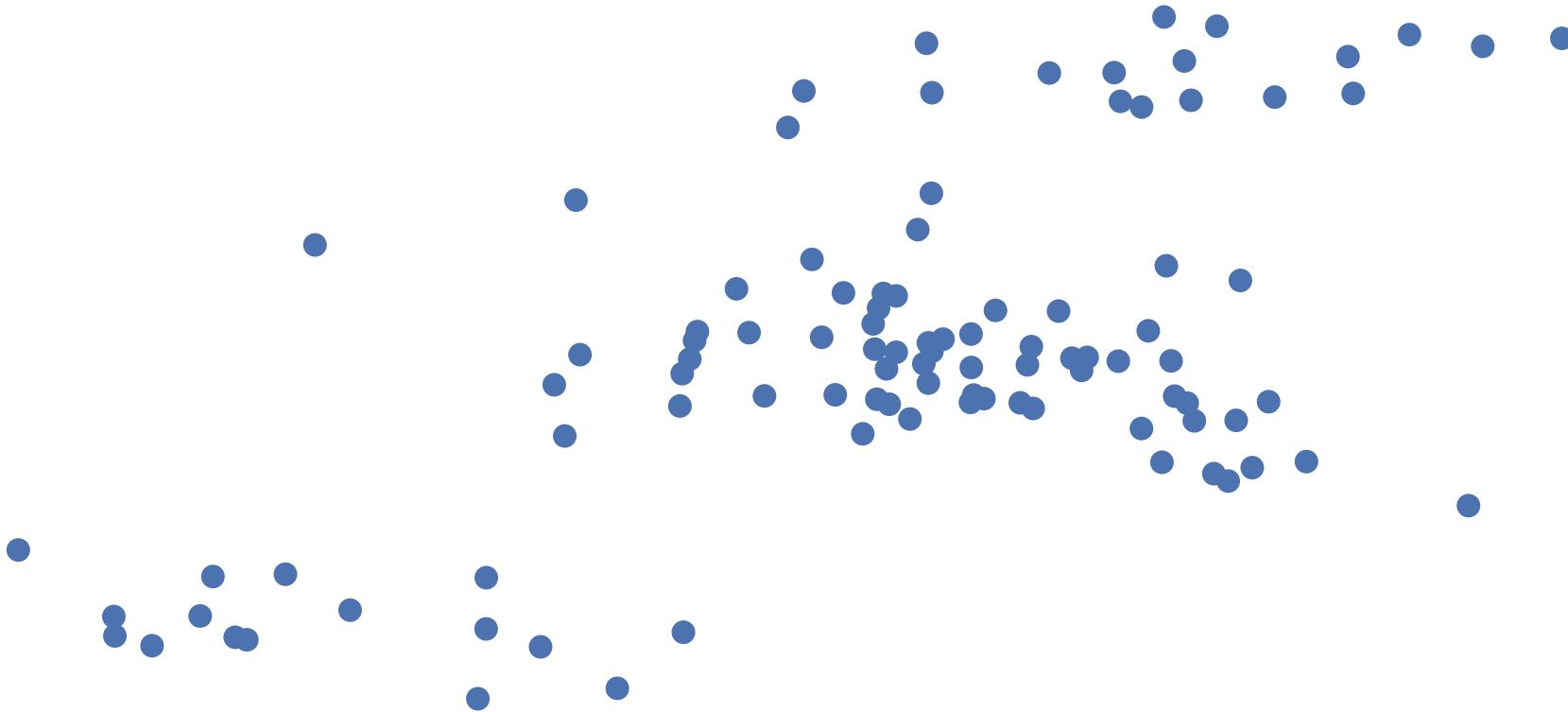
Is artificial intelligence set to become
art's next medium?

AI 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

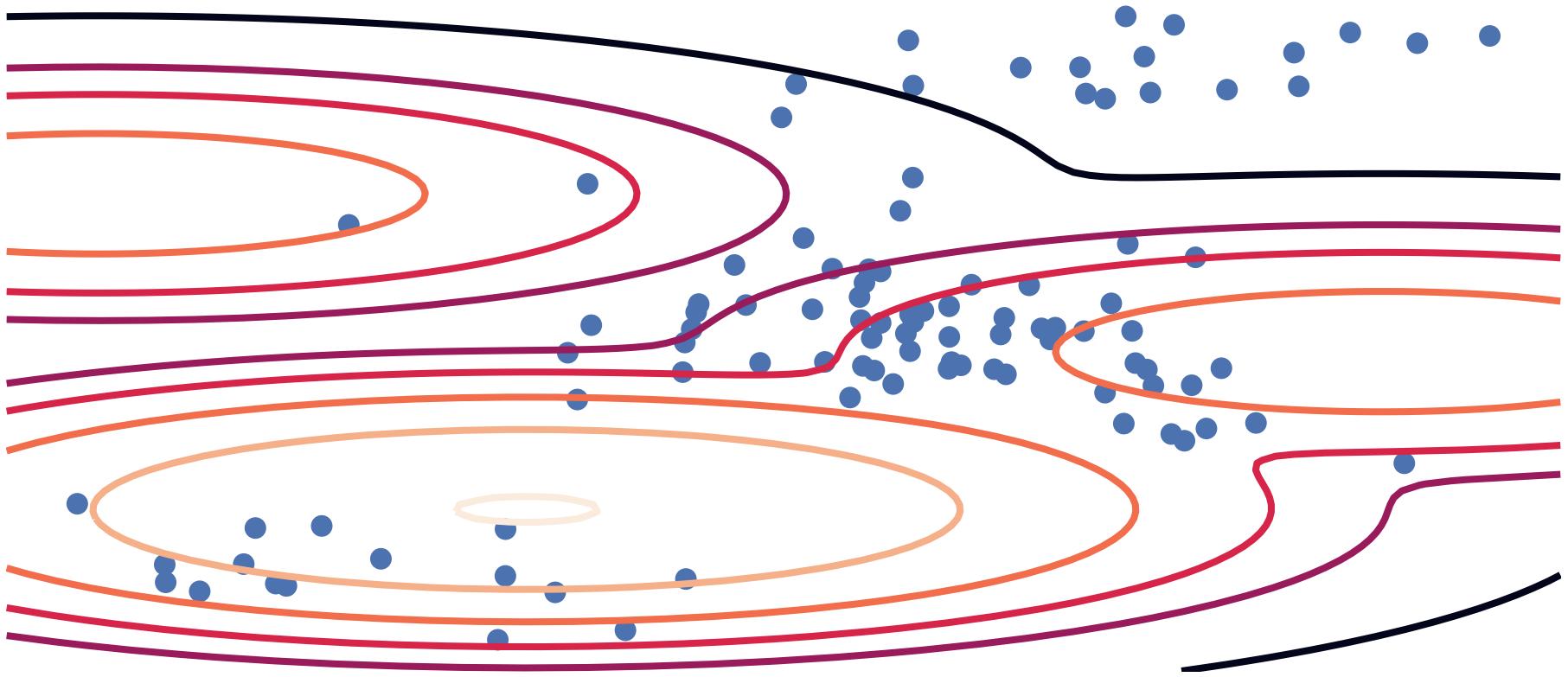
The portrait in its gilt frame depicts a portly gentleman, possibly French and — to judge by his dark frockcoat and plain white collar — a man of the church. The work appears unfinished: the facial features are somewhat indistinct and there are blank areas of canvas. Oddly, the whole composition is displaced 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:

$$\min_G \max_D \mathbb{E}_x[\log(D(x))] + \mathbb{E}_z[\log(1 - D(G(z)))]$$

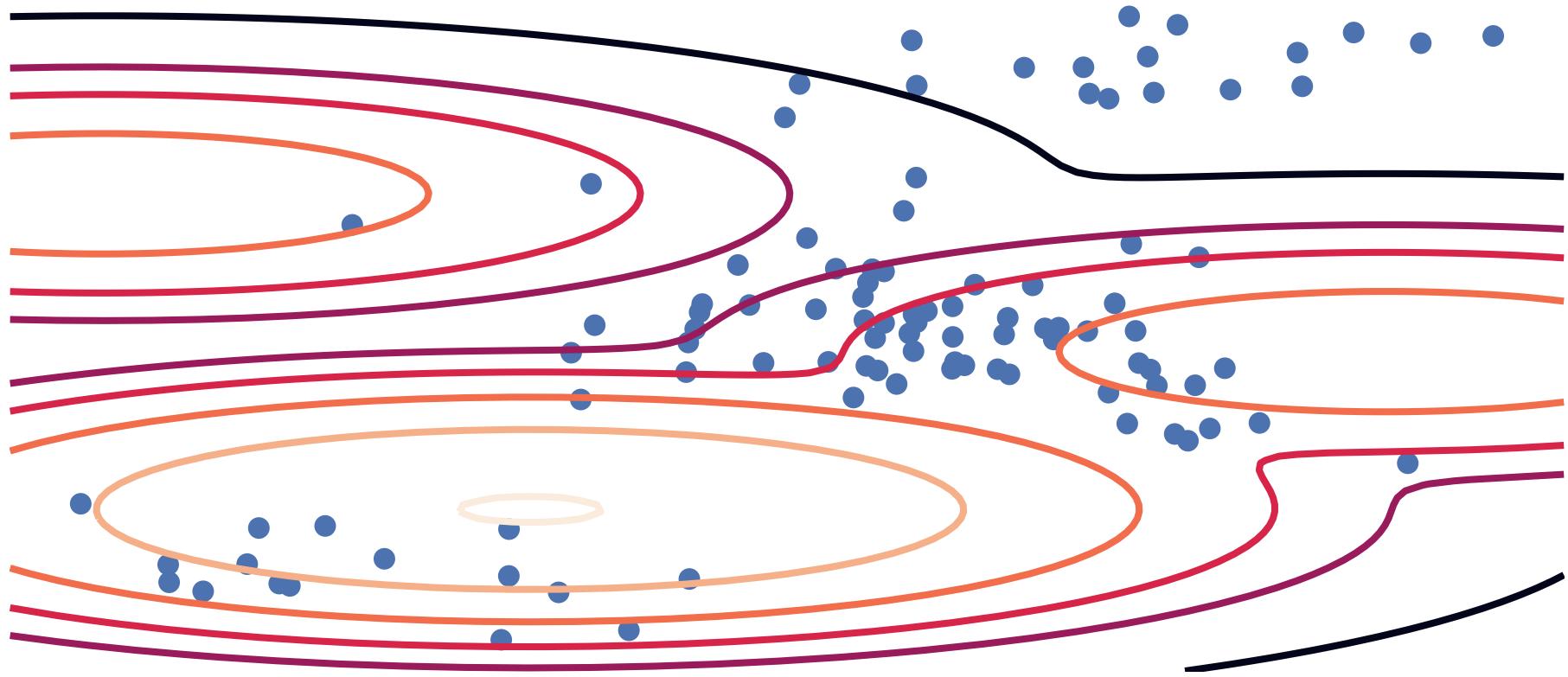
Generative models: a traditional way



Generative models: a traditional way

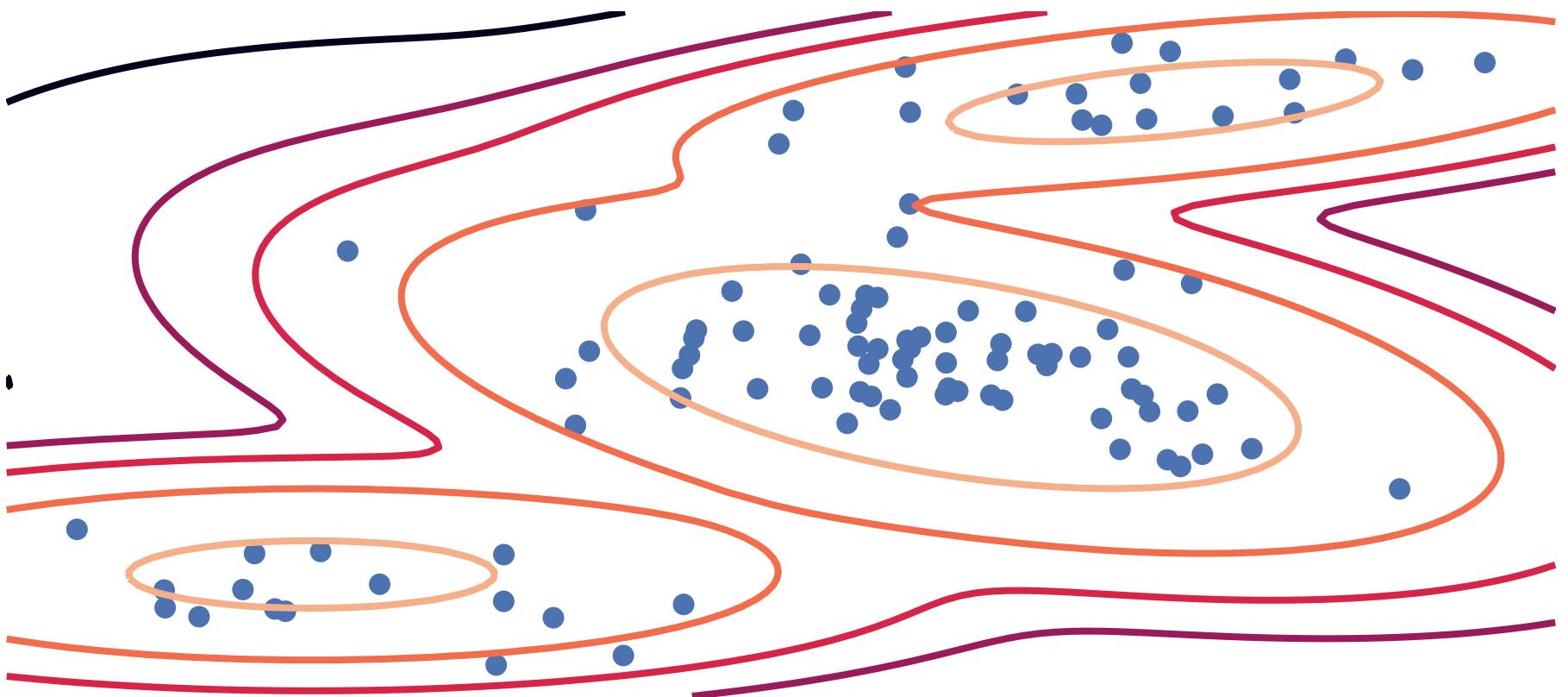


Generative models: a traditional way



Maximum likelihood: $\max_{\theta} \mathbb{E}_{X \sim P} [\log q_{\theta}(X)]$

Generative models: a traditional way



Maximum likelihood: $\max_{\theta} \mathbb{E}_{X \sim P} [\log q_{\theta}(X)]$

Kullback-Liebler (KL) divergence

- $\text{KL}(\mathbb{P}\|\mathbb{Q}) = \int p(x) \log \frac{p(x)}{q(x)} dx$
- A measure of how different two distributions are
- Related to entropy, etc: average number of bits to encode a sample from \mathbb{P} with a code for \mathbb{Q}

Kullback-Liebler (KL) divergence

- $\text{KL}(\mathbb{P} \parallel \mathbb{Q}) = \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\mathbb{q}(x)} dx$
- A measure of how different two distributions are
- Related to entropy, etc: average number of bits to encode a sample from \mathbb{P} with a code for \mathbb{Q}
- Minimizing one direction gives you maximum likelihood:

$$\arg \min_{\theta} \text{KL}(\mathbb{P} \parallel \mathbb{Q}_{\theta}) = \arg \min_{\theta} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\mathbb{q}_{\theta}(x)} dx$$

Kullback-Liebler (KL) divergence

- $\text{KL}(\mathbb{P} \parallel \mathbb{Q}) = \int p(x) \log \frac{p(x)}{q(x)} dx$
- A measure of how different two distributions are
- Related to entropy, etc: average number of bits to encode a sample from \mathbb{P} with a code for \mathbb{Q}
- Minimizing one direction gives you maximum likelihood:

$$\begin{aligned}\arg \min_{\theta} \text{KL}(\mathbb{P} \parallel \mathbb{Q}_{\theta}) &= \arg \min_{\theta} \int p(x) \log \frac{p(x)}{q_{\theta}(x)} dx \\ &= \arg \min_{\theta} -H[\mathbb{P}] - \mathbb{E}_{X \sim \mathbb{P}} [\log q_{\theta}(X)]\end{aligned}$$

Kullback-Liebler (KL) divergence

- $\text{KL}(\mathbb{P} \parallel \mathbb{Q}) = \int p(x) \log \frac{p(x)}{q(x)} dx$
- A measure of how different two distributions are
- Related to entropy, etc: average number of bits to encode a sample from \mathbb{P} with a code for \mathbb{Q}
- Minimizing one direction gives you maximum likelihood:

$$\begin{aligned}\arg \min_{\theta} \text{KL}(\mathbb{P} \parallel \mathbb{Q}_{\theta}) &= \arg \min_{\theta} \int p(x) \log \frac{p(x)}{q_{\theta}(x)} dx \\ &= \arg \min_{\theta} -H[\mathbb{P}] - \mathbb{E}_{X \sim \mathbb{P}} [\log q_{\theta}(X)] \\ &= \arg \max_{\theta} \mathbb{E}_{X \sim \mathbb{P}} [\log q_{\theta}(X)]\end{aligned}$$

Traditional models for images

- 1987-style generative model of faces (Eigenface via [Alex Egg](#))



Traditional models for images

- 1987-style generative model of faces (Eigenface via [Alex Egg](#))



- Can do fancier versions, of course...

Traditional models for images

- 1987-style generative model of faces (Eigenface via [Alex Egg](#))



- Can do fancier versions, of course...
- Usually based on Gaussian noise $\approx L_2$ loss

A hard case for traditional approaches

- One use case of generative models is inpainting [[Harry Yang](#)]:

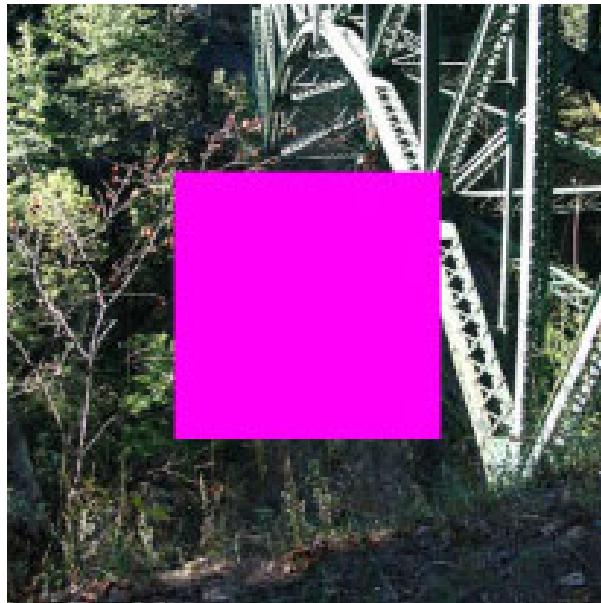


A hard case for traditional approaches

- One use case of generative models is inpainting [[Harry Yang](#)]:

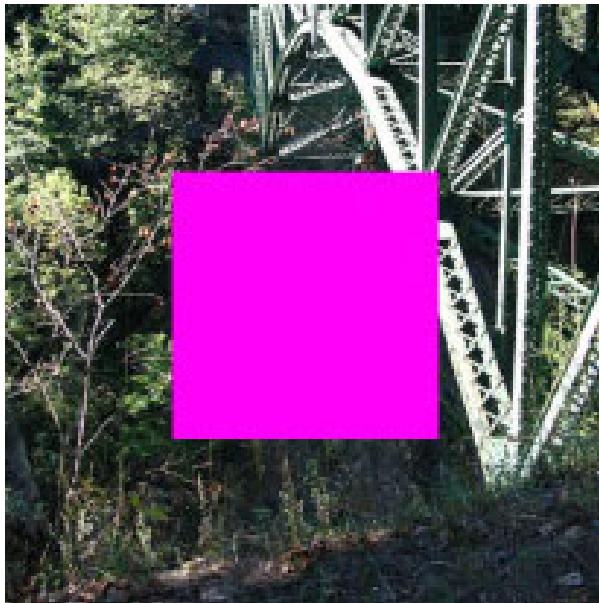
A hard case for traditional approaches

- One use case of generative models is inpainting [[Harry Yang](#)]:



A hard case for traditional approaches

- One use case of generative models is inpainting [[Harry Yang](#)]:



- L_2 loss / Gaussians will pick the *mean* of possibilities

A hard case for traditional approaches

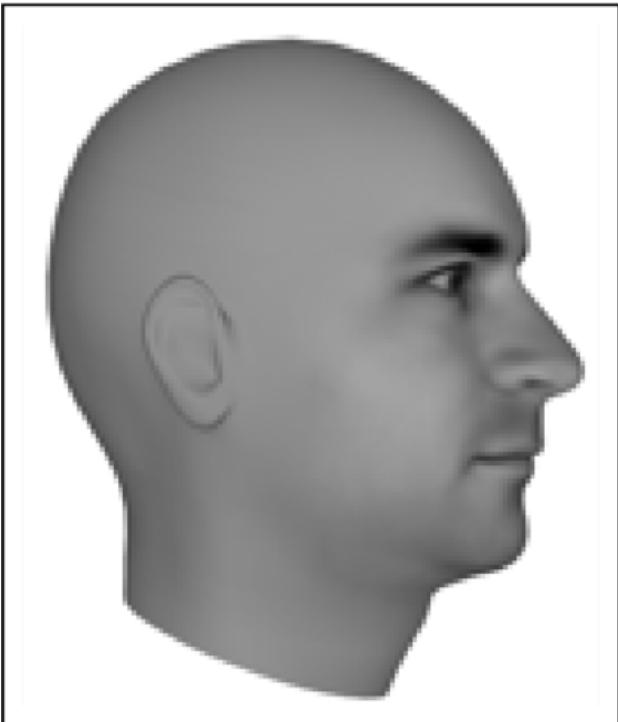
- One use case of generative models is inpainting [[Harry Yang](#)]:



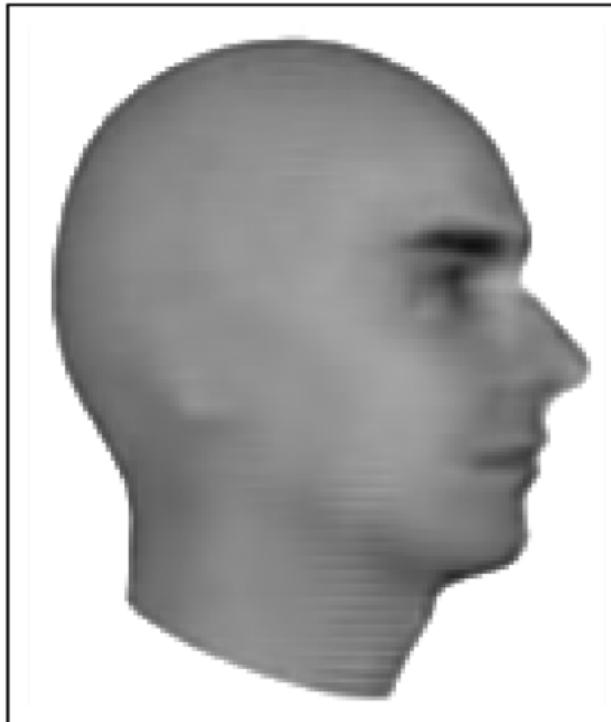
- L_2 loss / Gaussians will pick the *mean* of possibilities

Next-frame video prediction

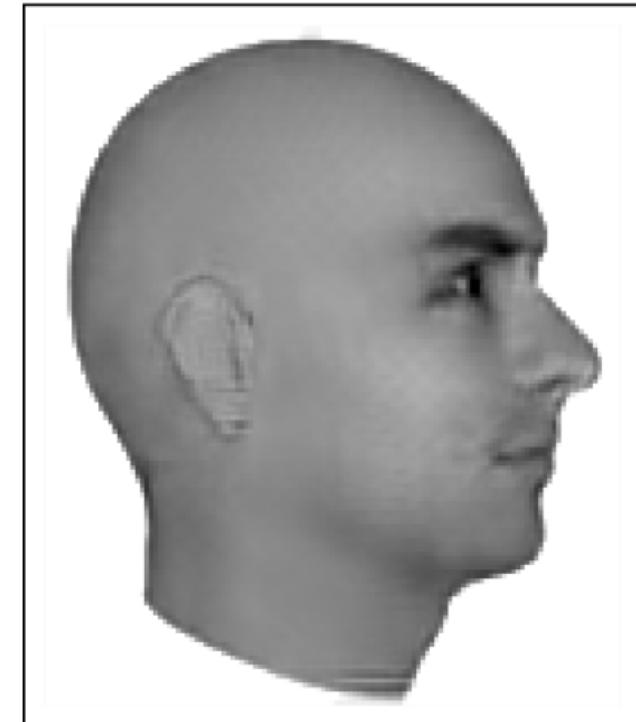
Ground Truth



MSE



Adversarial



[Lotter+ 2016]

Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (\mathbb{Q}_θ)



Discriminator



Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (\mathbb{Q}_θ)



Discriminator



Is this real?



Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?



Discriminator



Target (\mathbb{P})



Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?



Discriminator



Target (\mathbb{P})



No way! $\text{Pr}(\text{real}) = 0.03$

Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?
:(I'll try harder...

Discriminator



Target (\mathbb{P})



No way! $\Pr(\text{real}) = 0.03$

Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?



:(I'll try harder...

:

Discriminator



Target (\mathbb{P})



No way! $\text{Pr}(\text{real}) = 0.03$

Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?



:(
I'll try harder...

Is this real?



:

Discriminator



Target (\mathbb{P})



No way! $\text{Pr}(\text{real}) = 0.03$

Trick a discriminator [Goodfellow+ NeurIPS-14]

Generator (Q_θ)



Is this real?



:(
I'll try harder...

Is this real?



Discriminator



Target (\mathbb{P})

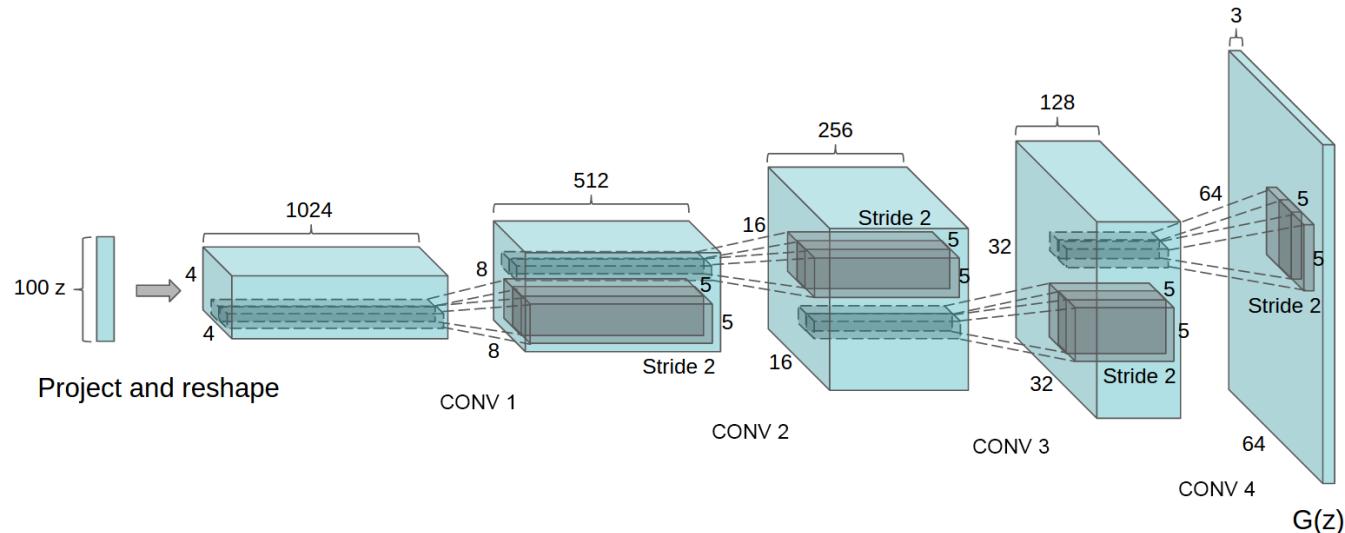
No way! $\text{Pr}(\text{real}) = 0.03$

:

Umm... $\text{Pr}(\text{real}) = 0.48$

Generator networks

- How to specify Q_θ ?



[Radford+ ICLR-16]

- $Z \sim \mathbb{Z} = \text{Uniform } ([-1, 1]^{100})$
- $G_\theta : [-1, 1]^{100} \rightarrow \mathcal{X}, G_\theta(Z) \sim Q_\theta$

GANs in equations

- Tricking the discriminator:

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

GANs in equations

- Tricking the discriminator:

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

- Using the generator network for Q_{θ} :

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Z \sim \mathbb{Z}} [\log(1 - D_{\psi}(G_{\theta}(Z)))]$$

GANs in equations

- Tricking the discriminator:

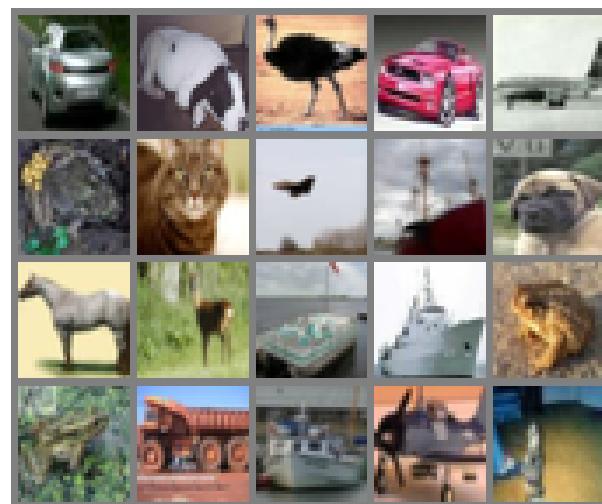
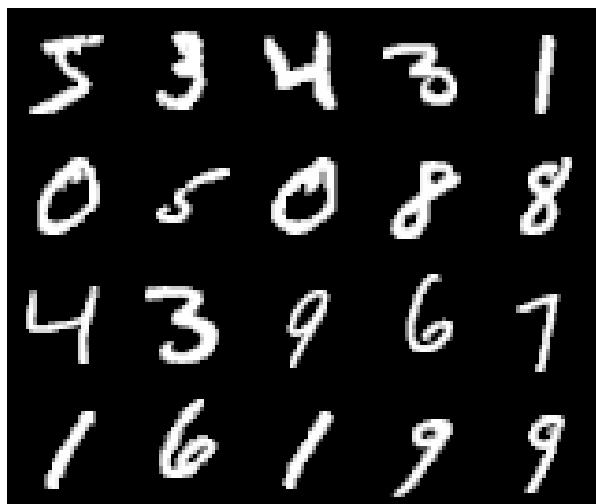
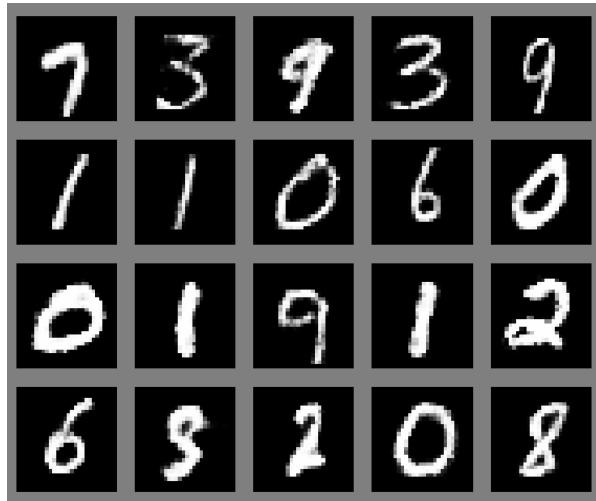
$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

- Using the generator network for Q_{θ} :

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Z \sim \mathbb{Z}} [\log(1 - D_{\psi}(G_{\theta}(Z)))]$$

- Can do alternating gradient descent!

Original paper's results [Goodfellow+ NeurIPS-14]

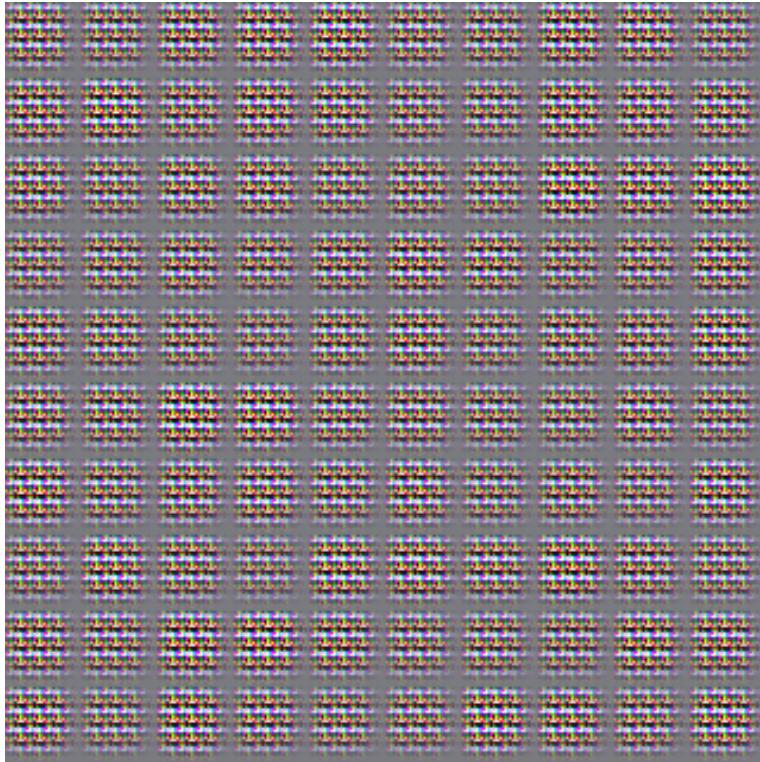


DCGAN results [Radford+ ICLR-16]



Training instability

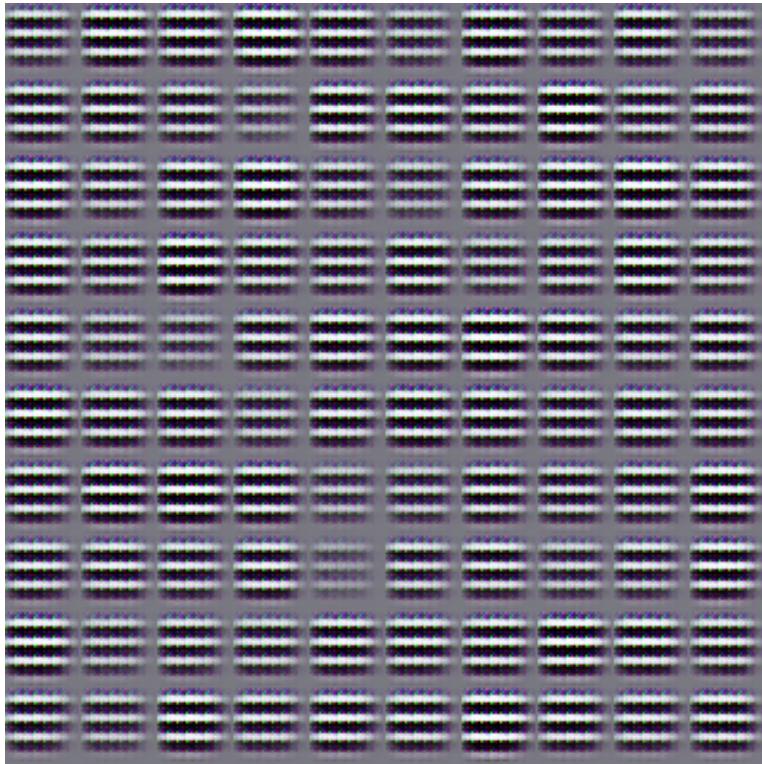
Running code from [[Salimans+ NeurIPS-16](#)]:



Run 1, epoch 1

Training instability

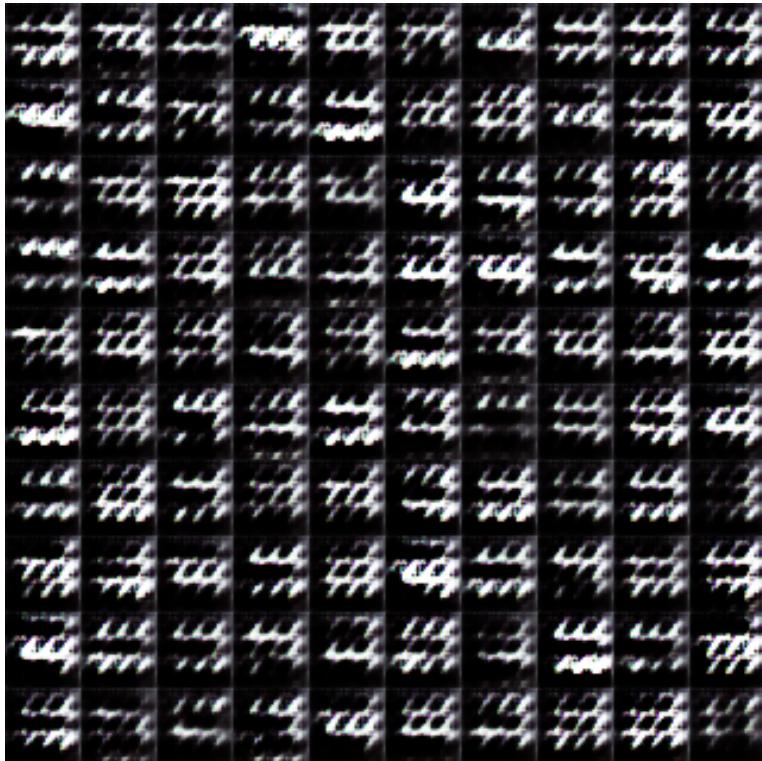
Running code from [[Salimans+ NeurIPS-16](#)]:



Run 1, epoch 2

Training instability

Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 3

Training instability

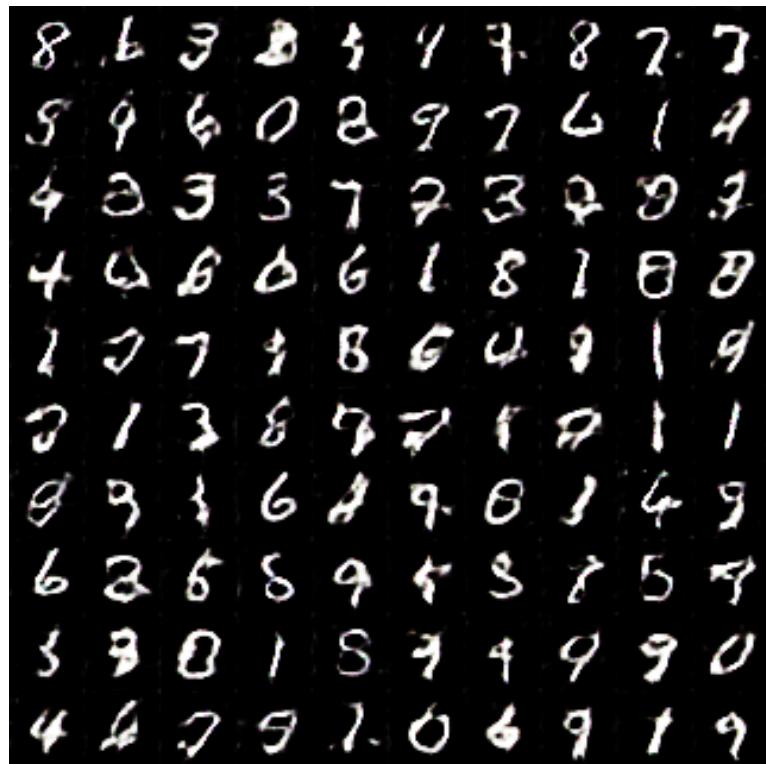
Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 4

Training instability

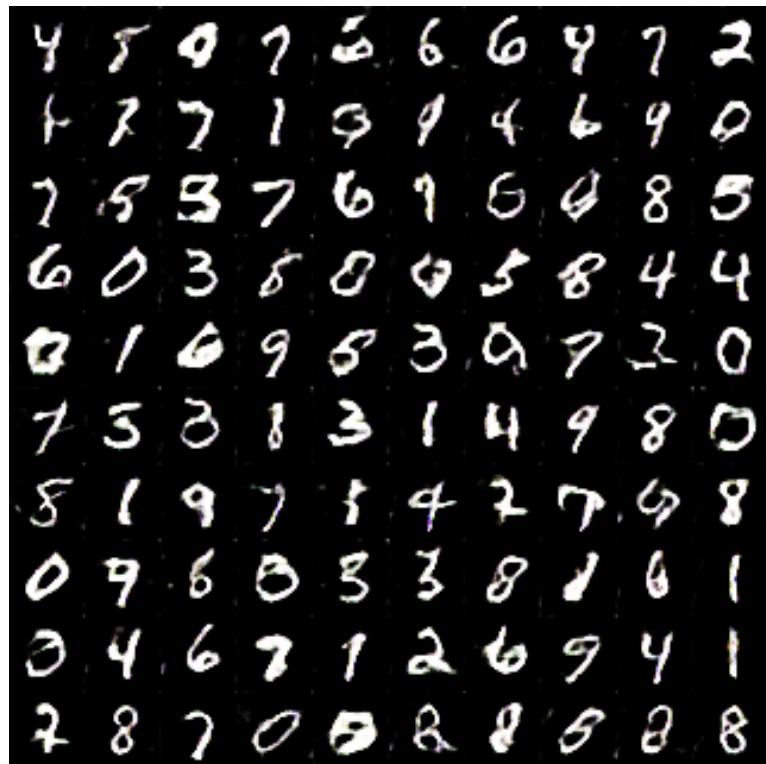
Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 5

Training instability

Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 6

Training instability

Running code from [Salimans+ NeurIPS-16]:

1	4	1	9	0	0	7	7	1	6
8	3	1	9	7	8	4	8	5	5
3	0	5	5	3	7	9	1	1	8
9	8	8	4	9	5	4	8	2	9
1	7	6	7	3	3	0	4	7	7
3	2	8	7	8	6	3	6	8	3
2	6	7	4	0	6	3	6	1	5
0	6	2	2	9	2	4	6	1	8
2	3	5	6	7	4	7	4	3	1
9	9	6	4	8	9	8	4	5	6

Run 1, epoch 11

Training instability

Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 501

Training instability

Running code from [Salimans+ NeurIPS-16]:



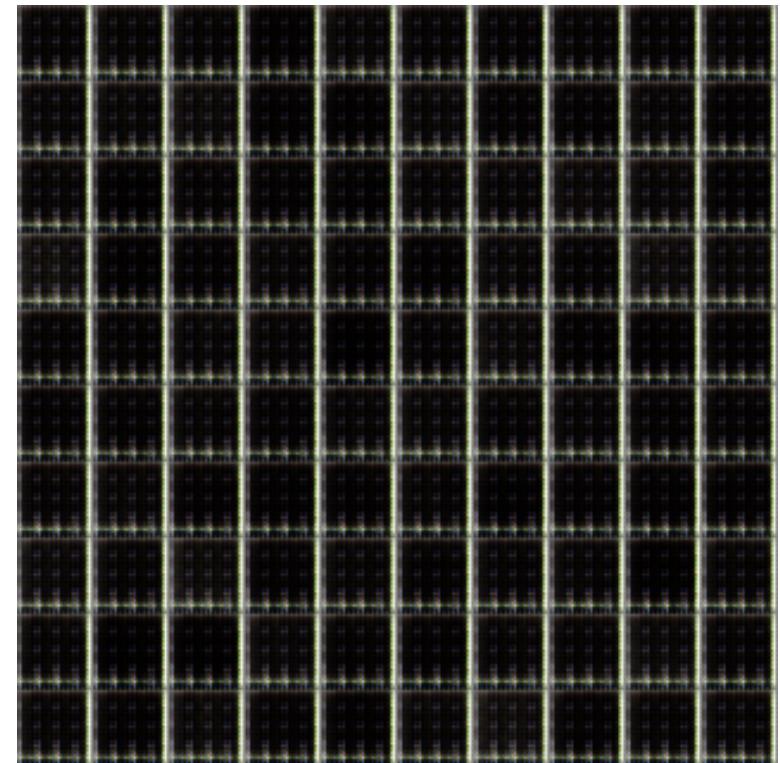
Run 1, epoch 900

Training instability

Running code from [[Salimans+ NeurIPS-16](#)]:



Run 1, epoch 900



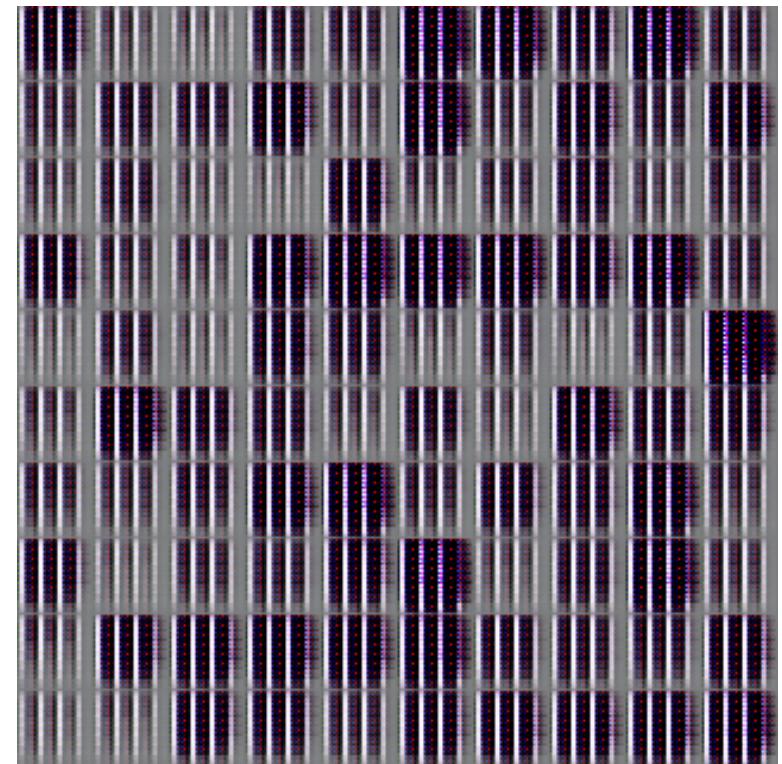
Run 2, epoch 1

Training instability

Running code from [[Salimans+ NeurIPS-16](#)]:



Run 1, epoch 900



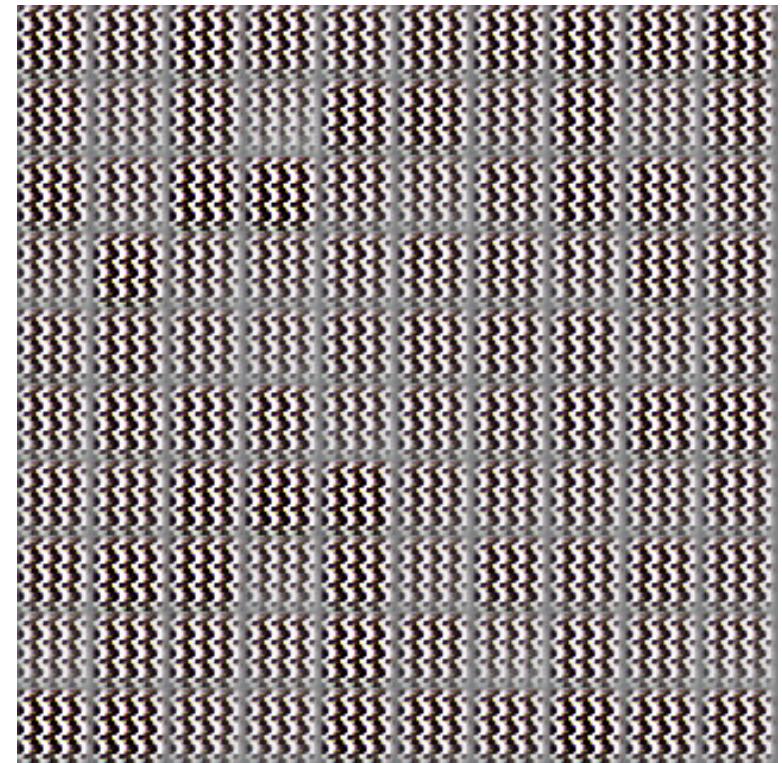
Run 2, epoch 2

Training instability

Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 900



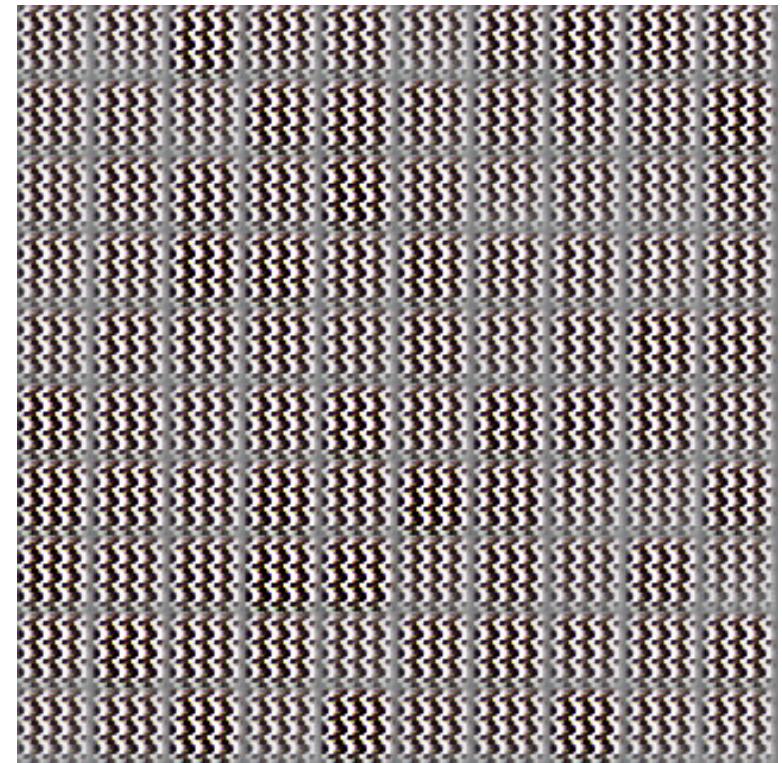
Run 2, epoch 3

Training instability

Running code from [Salimans+ NeurIPS-16]:



Run 1, epoch 900



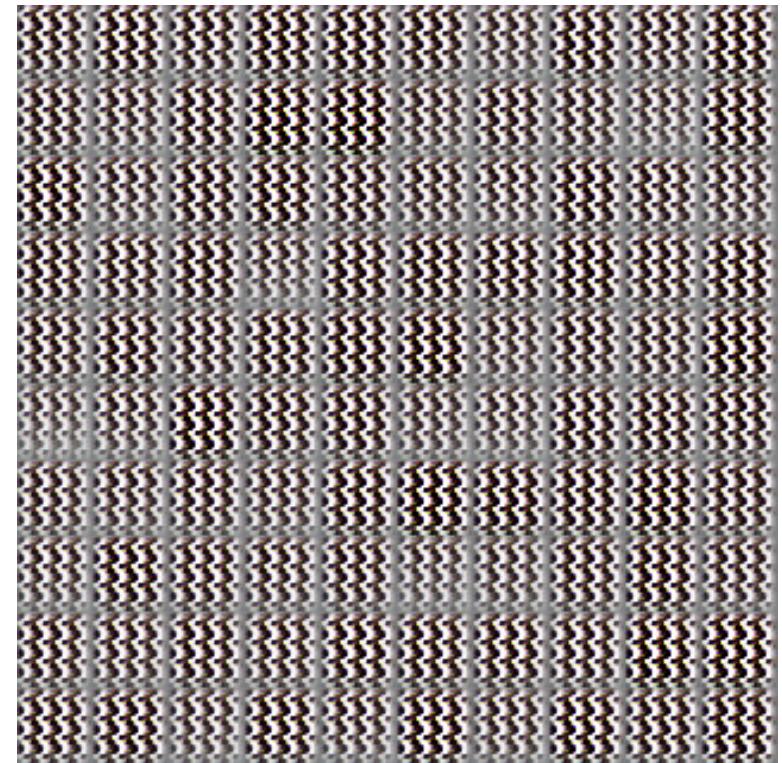
Run 2, epoch 4

Training instability

Running code from [[Salimans+ NeurIPS-16](#)]:



Run 1, epoch 900



Run 2, epoch 5

One view: distances between distributions

- What happens when D_ψ is at its optimum?

One view: distances between distributions

- What happens when D_{ψ} is at its optimum?
- If distributions have densities, $D_{\psi}^*(x) = \frac{\textcolor{blue}{p}(x)}{\textcolor{blue}{p}(x) + \textcolor{brown}{q}_{\theta}(x)}$

One view: distances between distributions

- What happens when D_ψ is at its optimum?
- If distributions have densities, $D_\psi^*(x) = \frac{p(x)}{p(x) + q_\theta(x)}$
- If D_ψ stays optimal throughout, θ tries to minimize

$$\frac{1}{2} \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} \left[\log \frac{p(\mathbf{X})}{p(\mathbf{X}) + q_\theta(\mathbf{X})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y} \sim \mathbb{Q}_\theta} \left[\log \frac{q_\theta(\mathbf{X})}{p(\mathbf{X}) + q_\theta(\mathbf{X})} \right]$$

which is $\text{JS}(\mathbb{P}, \mathbb{Q}_\theta) - \log 2$

Jensen-Shannon divergence

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbf{\color{blue} p}(x) \log \frac{\mathbf{\color{blue} p}(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbf{\color{brown} q}_\theta(x) \log \frac{\mathbf{\color{brown} q}_\theta(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \end{aligned}$$

Jensen-Shannon divergence

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbb{q}_\theta(x) \log \frac{\mathbb{q}_\theta(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &= \frac{1}{2} \text{KL}\left(\mathbb{P} \parallel \frac{\mathbb{P} + \mathbb{Q}_\theta}{2}\right) + \frac{1}{2} \text{KL}\left(\mathbb{Q}_\theta \parallel \frac{\mathbb{P} + \mathbb{Q}_\theta}{2}\right) \end{aligned}$$

Jensen-Shannon divergence

$$\begin{aligned}\text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbb{q}_\theta(x) \log \frac{\mathbb{q}_\theta(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &= \frac{1}{2} \text{KL}\left(\mathbb{P} \parallel \frac{\mathbb{P} + \mathbb{Q}_\theta}{2}\right) + \frac{1}{2} \text{KL}\left(\mathbb{Q}_\theta \parallel \frac{\mathbb{P} + \mathbb{Q}_\theta}{2}\right) \\ &= \text{H}\left[\frac{\mathbb{P} + \mathbb{Q}_\theta}{2}\right] - \frac{\text{H}[\mathbb{P}] + \text{H}[\mathbb{Q}_\theta]}{2}\end{aligned}$$

JS with disjoint support [Arjovsky/Bottou ICLR-17]

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbf{\color{blue} p}(x) \log \frac{\mathbf{\color{blue} p}(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbf{\color{brown} q}_\theta(x) \log \frac{\mathbf{\color{brown} q}_\theta(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \end{aligned}$$

- If \mathbb{P} and \mathbb{Q}_θ have (almost) disjoint support

$$\frac{1}{2} \int \mathbf{\color{blue} p}(x) \log \frac{\mathbf{\color{blue} p}(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x)} dx$$

JS with disjoint support [Arjovsky/Bottou ICLR-17]

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbb{q}_\theta(x) \log \frac{\mathbb{q}_\theta(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \end{aligned}$$

- If \mathbb{P} and \mathbb{Q}_θ have (almost) disjoint support

$$\frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x)} dx = \frac{1}{2} \int \mathbb{p}(x) \log(2) dx$$

JS with disjoint support [Arjovsky/Bottou ICLR-17]

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbb{q}_\theta(x) \log \frac{\mathbb{q}_\theta(x)}{\frac{1}{2}\mathbb{p}(x) + \frac{1}{2}\mathbb{q}_\theta(x)} dx \end{aligned}$$

- If \mathbb{P} and \mathbb{Q}_θ have (almost) disjoint support

$$\frac{1}{2} \int \mathbb{p}(x) \log \frac{\mathbb{p}(x)}{\frac{1}{2}\mathbb{p}(x)} dx = \frac{1}{2} \int \mathbb{p}(x) \log(2) dx = \frac{1}{2} \log 2$$

JS with disjoint support [Arjovsky/Bottou ICLR-17]

$$\begin{aligned} \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) &= \frac{1}{2} \int \mathbf{\color{blue} p}(x) \log \frac{\mathbf{\color{blue} p}(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \\ &\quad + \frac{1}{2} \int \mathbf{\color{brown} q}_\theta(x) \log \frac{\mathbf{\color{brown} q}_\theta(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x) + \frac{1}{2}\mathbf{\color{brown} q}_\theta(x)} dx \end{aligned}$$

- If \mathbb{P} and \mathbb{Q}_θ have (almost) disjoint support

$$\frac{1}{2} \int \mathbf{\color{blue} p}(x) \log \frac{\mathbf{\color{blue} p}(x)}{\frac{1}{2}\mathbf{\color{blue} p}(x)} dx = \frac{1}{2} \int \mathbf{\color{blue} p}(x) \log(2) dx = \frac{1}{2} \log 2$$

$$\text{so } \text{JS}(\mathbb{P}, \mathbb{Q}_\theta) = \log 2$$

Discriminator point of view

Generator (Q_θ)



Discriminator



Discriminator point of view

Generator (Q_θ)



Discriminator



Is this real?



Discriminator point of view

Generator (Q_θ)



Is this real?



Discriminator



Target (\mathbb{P})



Discriminator point of view

Generator (Q_θ)



Is this real?



Discriminator



Target (\mathbb{P})



No way! $\text{Pr}(\text{real}) = 0.00$

Discriminator point of view

Generator (Q_θ)



Is this real?



:(| I don't know how to do any better...

Discriminator



Target (\mathbb{P})



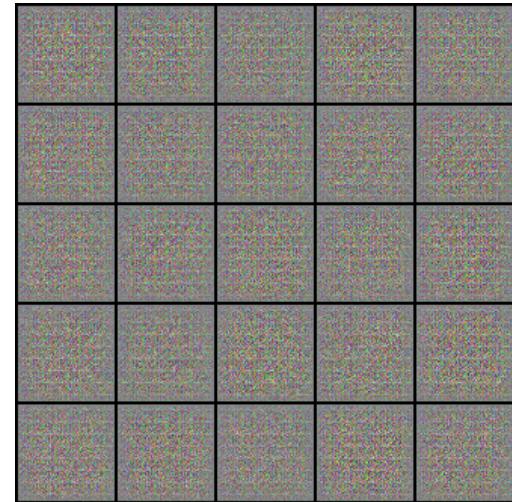
No way! $\text{Pr}(\text{real}) = 0.00$

How likely is disjoint support?

- At initialization, pretty reasonable:



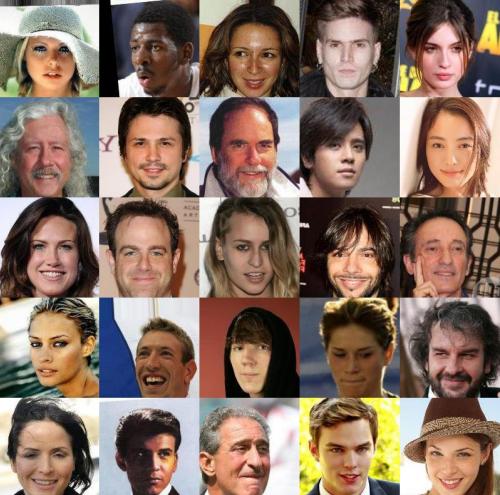
Q_θ :



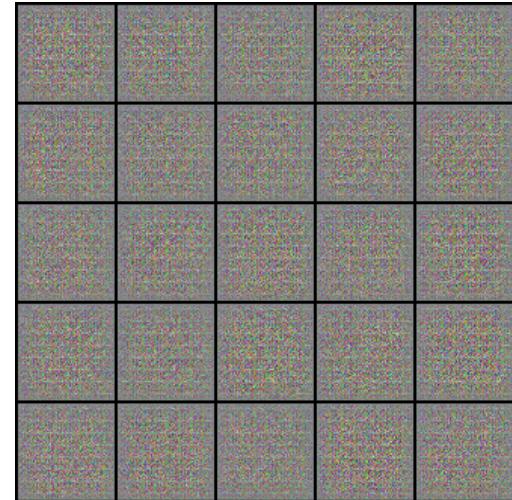
How likely is disjoint support?

- At initialization, pretty reasonable:

$P:$



$Q_\theta:$



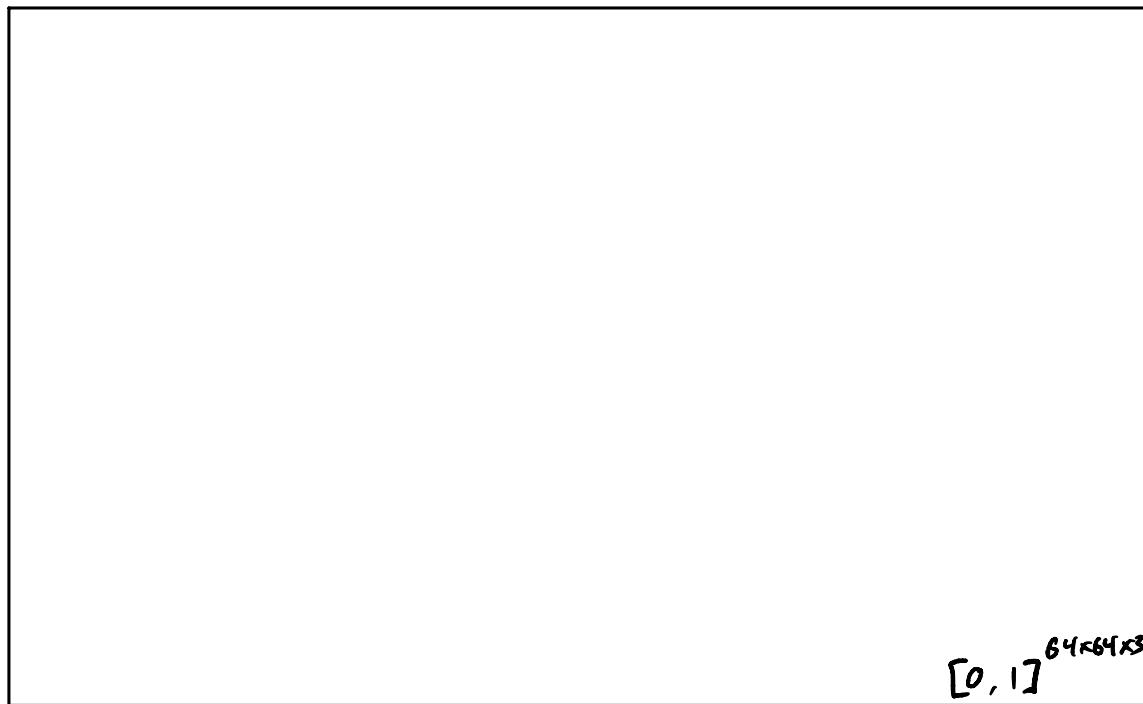
- Note that intersecting *almost nowhere* (with probability 0) doesn't change matters

Manifold structure

- If $G_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{64 \times 64 \times 3}$, for usual G_θ, Q_θ is supported on a countable union of manifolds with $\dim \leq 100$

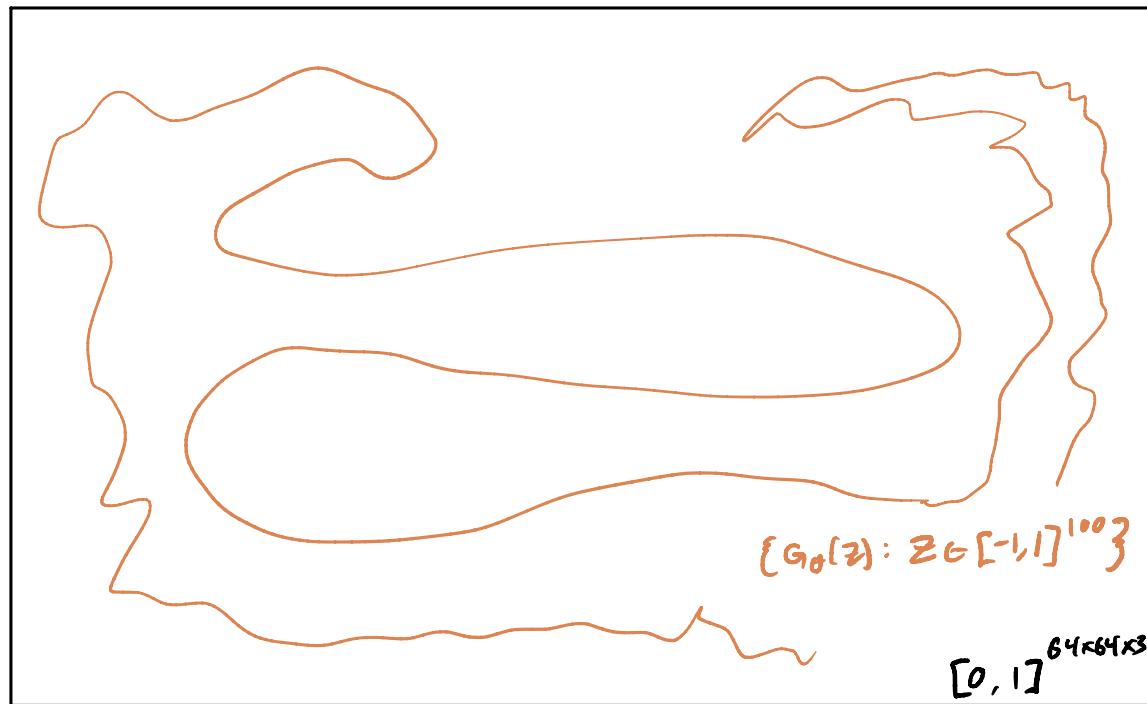
Manifold structure

- If $G_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{64 \times 64 \times 3}$, for usual G_θ, Q_θ is supported on a countable union of manifolds with $\dim \leq 100$



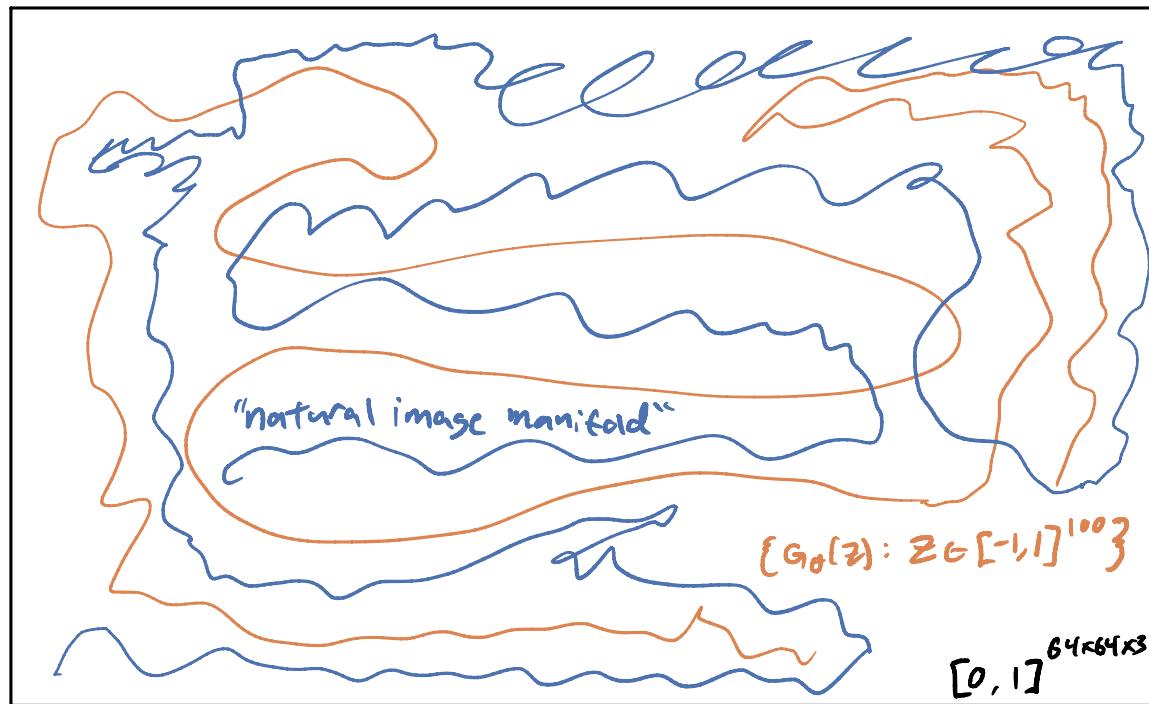
Manifold structure

- If $G_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{64 \times 64 \times 3}$, for usual G_θ , Q_θ is supported on a countable union of manifolds with $\dim \leq 100$



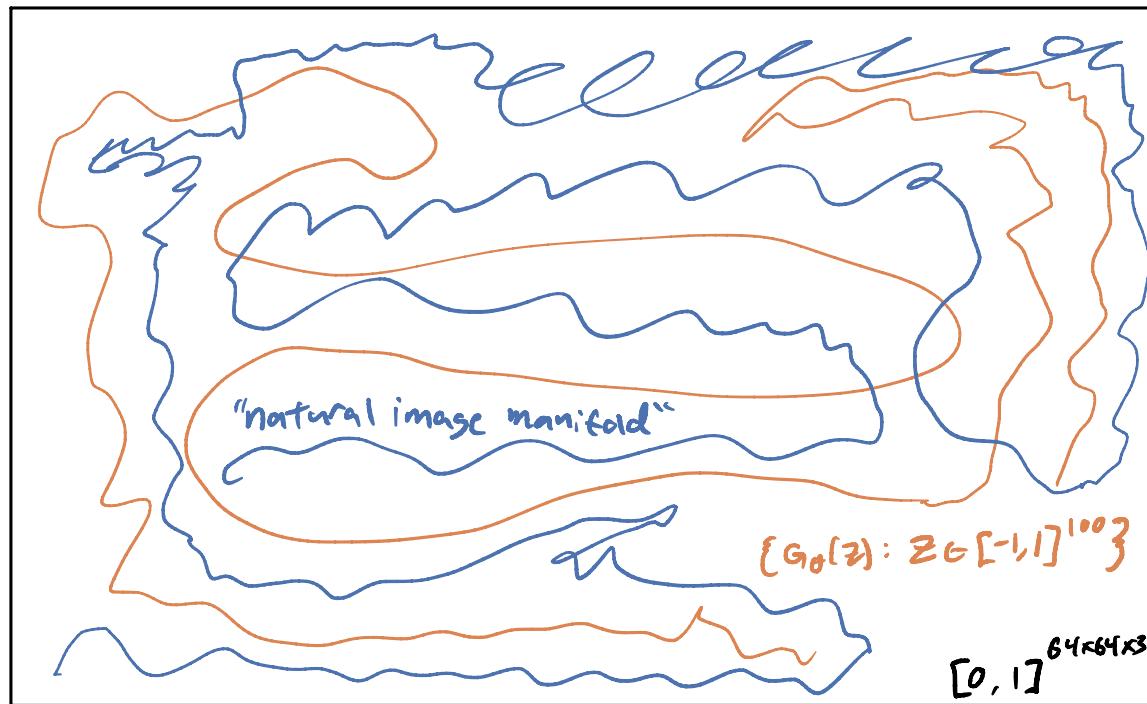
Manifold structure

- If $G_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{64 \times 64 \times 3}$, for usual G_θ , Q_θ is supported on a countable union of manifolds with $\dim \leq 100$
- “Natural image manifold” for \mathbb{P} usually considered low-dim



Manifold structure

- If $G_\theta : \mathbb{R}^{100} \rightarrow \mathbb{R}^{64 \times 64 \times 3}$, for usual G_θ , Q_θ is supported on a countable union of manifolds with $\dim \leq 100$
- “Natural image manifold” for P usually considered low-dim
- No chance that they'd align at init, so $\text{JS}(P, Q_\theta) = \log 2$



A heuristic partial workaround

- Original GANs almost never use the minimax game

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

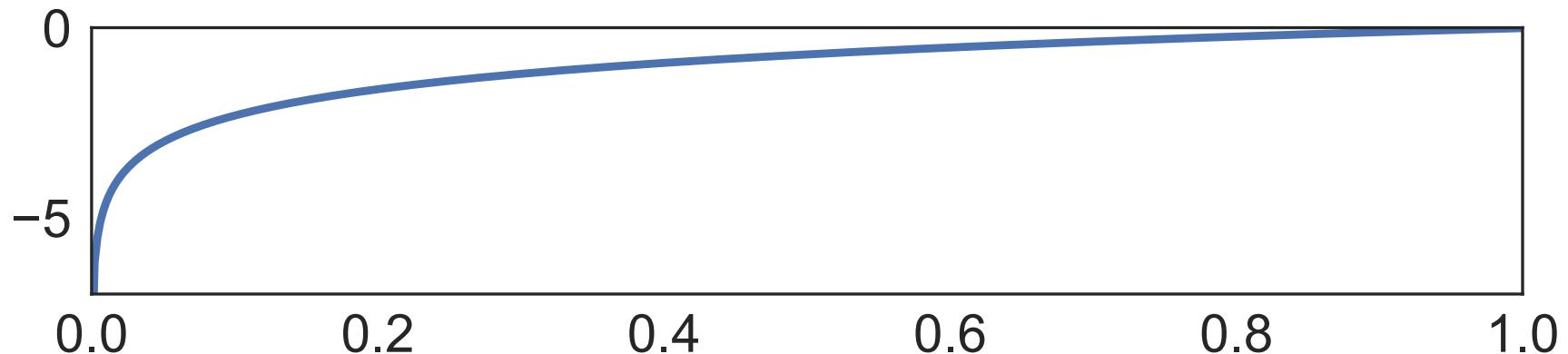
- $\max_{\theta} \log D_{\psi}(G_{\theta}(Z))$, not $\min_{\theta} \log(1 - D_{\psi}(G_{\theta}(Z)))$

A heuristic partial workaround

- Original GANs almost never use the minimax game

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

- $\max_{\theta} \log D_{\psi}(G_{\theta}(Z))$, not $\min_{\theta} \log(1 - D_{\psi}(G_{\theta}(Z)))$
- If D_{ψ} is near-perfect, near $\log 0$ instead of $\log 1$

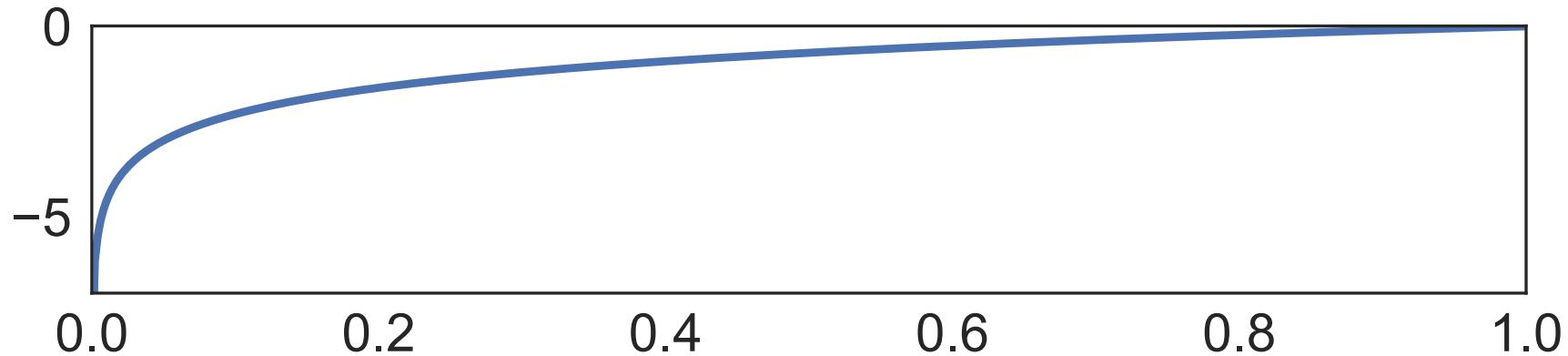


A heuristic partial workaround

- Original GANs almost never use the minimax game

$$\min_{\theta} \max_{\psi} \frac{1}{2} \mathbb{E}_{X \sim \mathbb{P}} [\log D_{\psi}(X)] + \frac{1}{2} \mathbb{E}_{Y \sim Q_{\theta}} [\log(1 - D_{\psi}(Y))]$$

- $\max_{\theta} \log D_{\psi}(G_{\theta}(Z))$, not $\min_{\theta} \log(1 - D_{\psi}(G_{\theta}(Z)))$
- If D_{ψ} is near-perfect, near $\log 0$ instead of $\log 1$



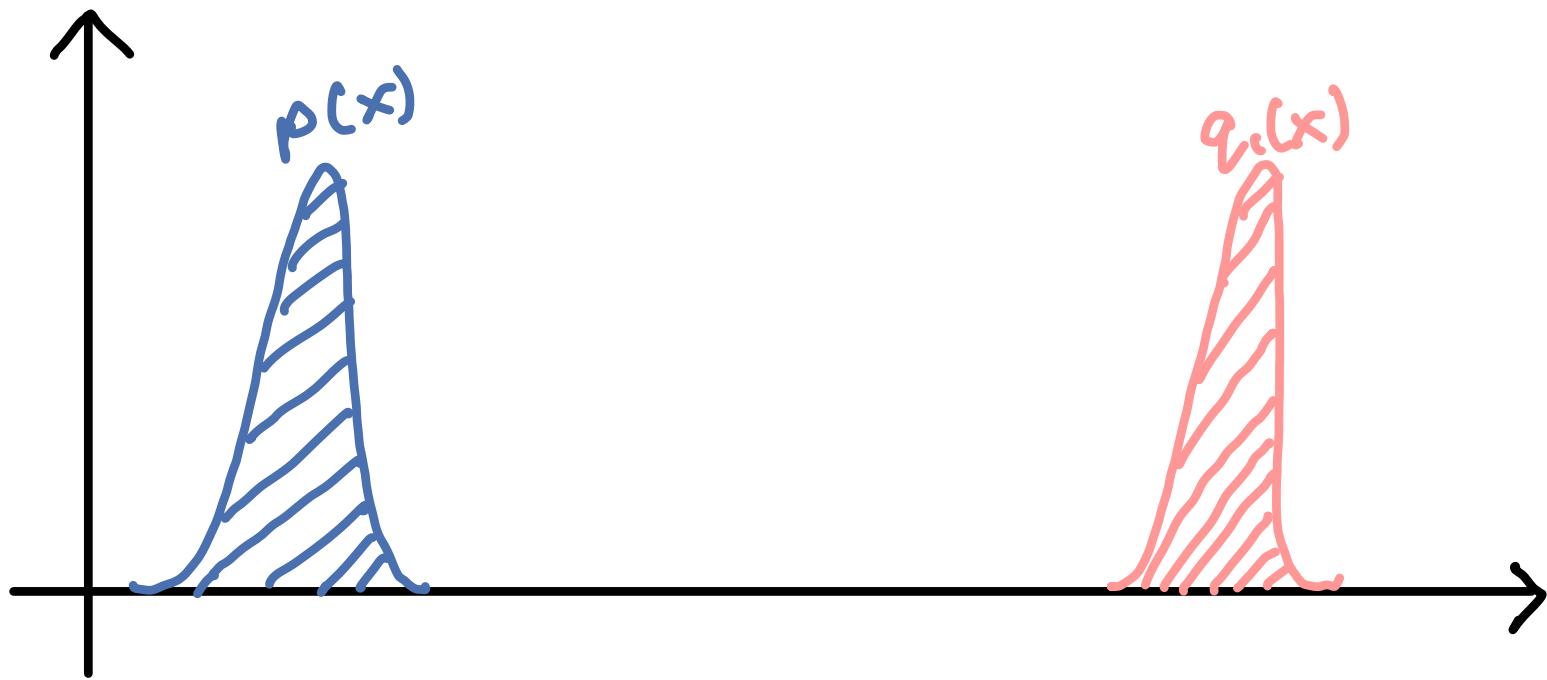
- When D_{ψ} is near-perfect, makes it unstable instead of stuck

Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately

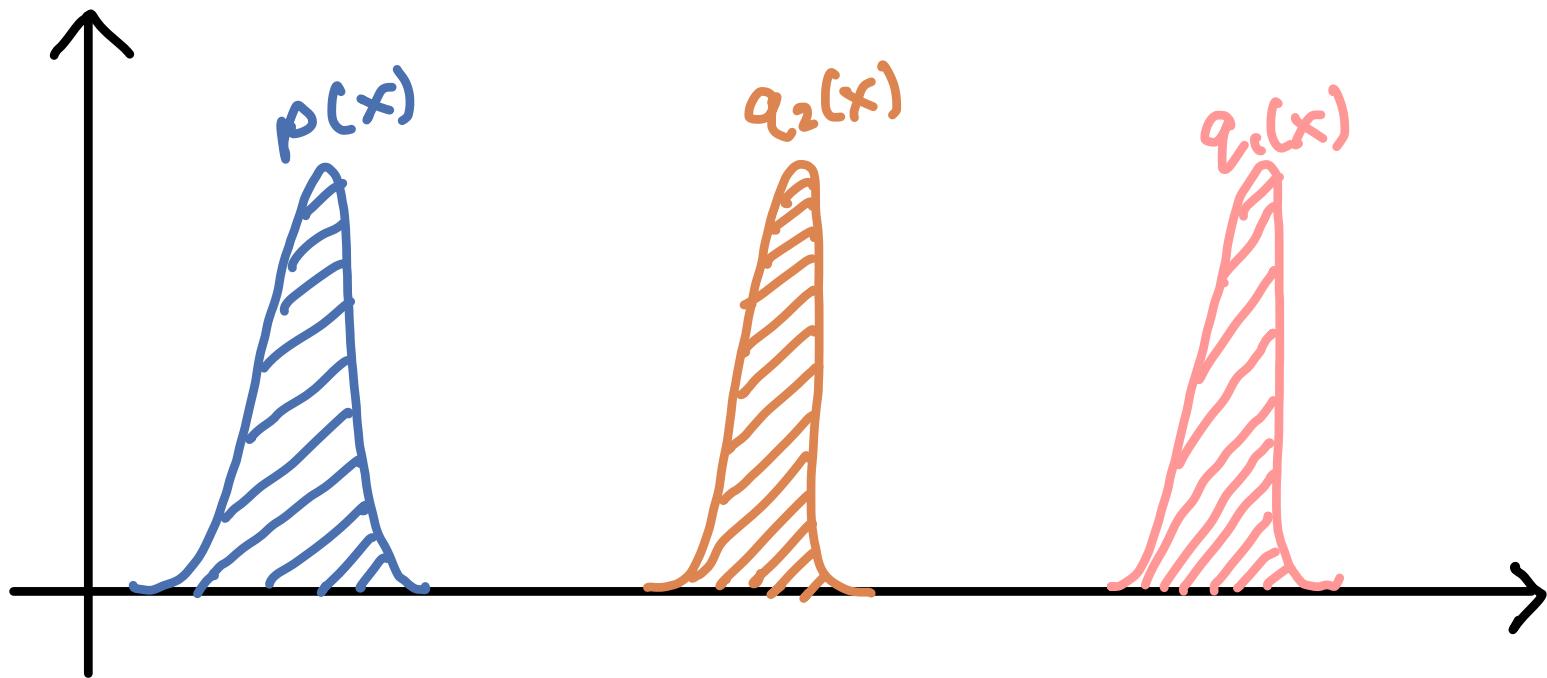
Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately



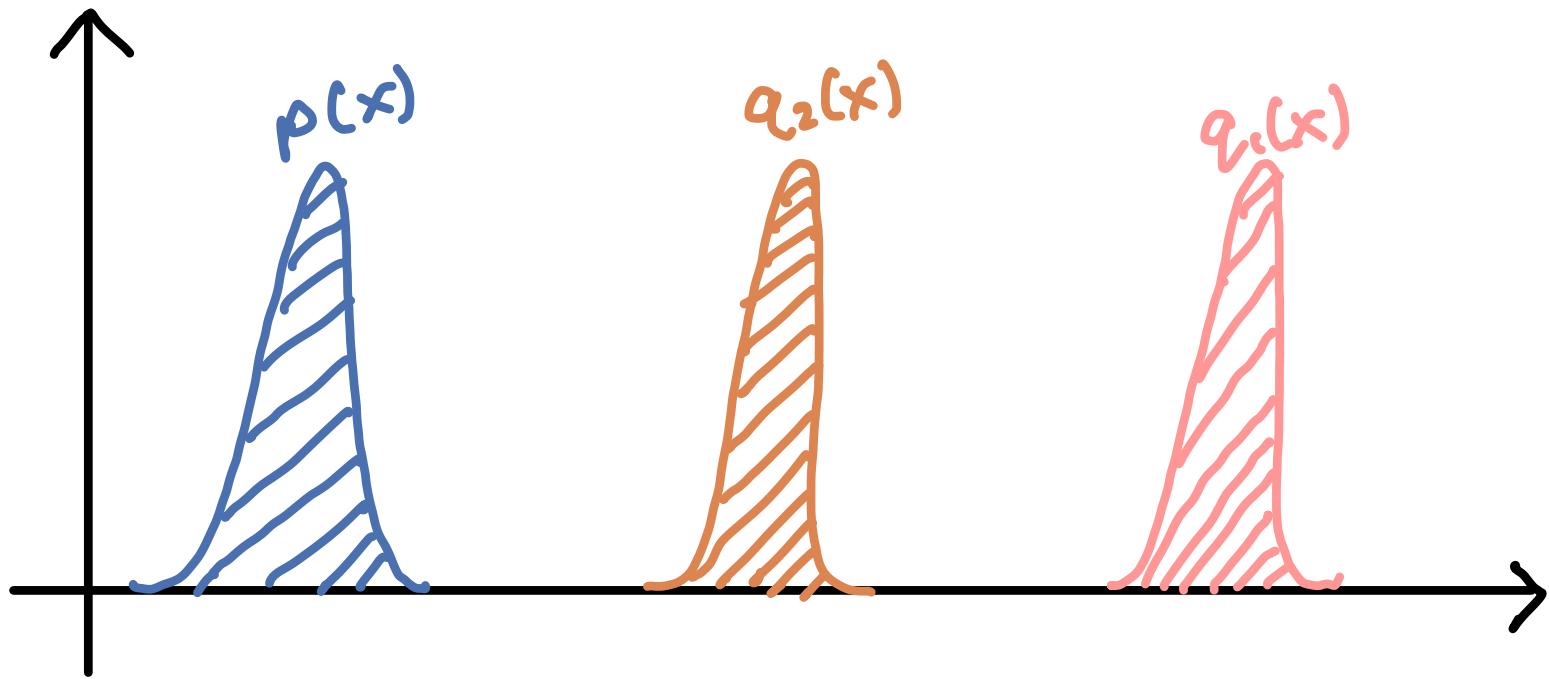
Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately



Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately
 - Can have $Q_n \rightarrow Q$ but $JS(Q_n, Q) = \log 2 \not\rightarrow 0$



Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately
 - Can have $Q_n \rightarrow Q$ but $JS(Q_n, Q) = \log 2 \not\rightarrow 0$
- Need some notion of geometry



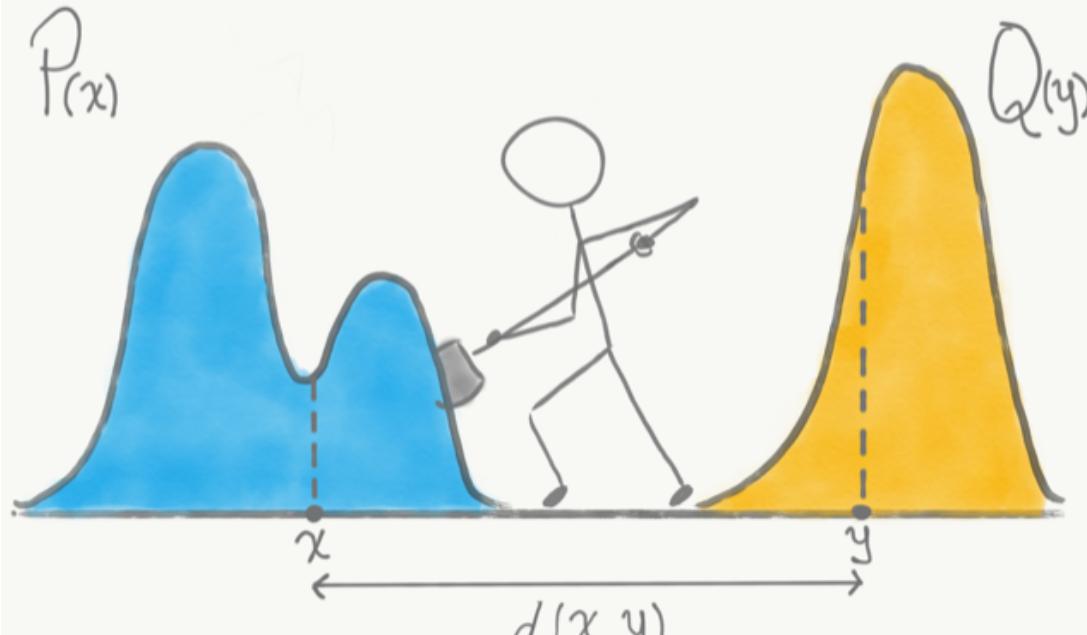
Real problem: KL/JS don't know geometry!

- $\int p(x) \log \frac{p(x)}{q(x)} dx$ treats each x totally separately
 - Can have $Q_n \rightarrow Q$ but $JS(Q_n, Q) = \log 2 \not\rightarrow 0$
- Need some notion of geometry
 - Should be *continuous in the weak topology*



Better Solution 1: Optimal Transport

Wasserstein, aka earth mover's distance:
How far do I have to move probability mass to turn \mathbb{P} into \mathbb{Q} ?



(from David Alvarez-Melis / Nicolo Fusi)

$$\mathcal{W}(\mathbb{P}, \mathbb{Q}) = \inf_C \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim C(\mathbb{P}, \mathbb{Q})} d(\mathbf{X}, \mathbf{Y})$$

(C is a joint distribution on (\mathbf{X}, \mathbf{Y}) with marginals \mathbb{P} and \mathbb{Q})

Wasserstein via Kantorovich-Rubinstein duality

$$\mathcal{W}(\mathbb{P}, \mathbb{Q}) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} [f(\mathbf{X})] - \mathbb{E}_{\mathbf{Y} \sim \mathbb{Q}} [f(\mathbf{Y})]$$

$f : \mathcal{X} \rightarrow \mathbb{R}$ is a 1-Lipschitz *critic function*

$$\|f\|_{\text{Lip}} = \sup_{x,y \in \mathcal{X}} \frac{|f(x) - f(y)|}{\|x - y\|} = \sup_{x \in \mathcal{X}} \|\nabla f(x)\|$$

\mathcal{W} is continuous (bounded domain): if $\mathbb{Q}_\theta \rightarrow \mathbb{P}$, then $\mathcal{W}(\mathbb{Q}_\theta, \mathbb{P}) \rightarrow 0$

Wasserstein via Kantorovich-Rubinstein duality

$$\mathcal{W}(\mathbb{P}, \mathbb{Q}) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} [f(\mathbf{X})] - \mathbb{E}_{\mathbf{Y} \sim \mathbb{Q}} [f(\mathbf{Y})]$$

$f : \mathcal{X} \rightarrow \mathbb{R}$ is a 1-Lipschitz *critic function*

$$\|f\|_{\text{Lip}} = \sup_{x,y \in \mathcal{X}} \frac{|f(x) - f(y)|}{\|x - y\|} = \sup_{x \in \mathcal{X}} \|\nabla f(x)\|$$



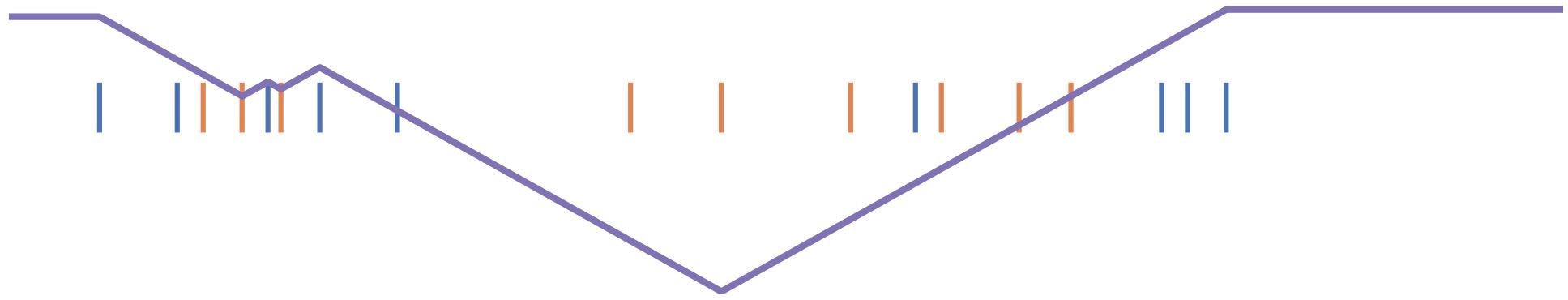
\mathcal{W} is continuous (bounded domain): if $\mathbb{Q}_\theta \rightarrow \mathbb{P}$, then $\mathcal{W}(\mathbb{Q}_\theta, \mathbb{P}) \rightarrow 0$

Wasserstein via Kantorovich-Rubinstein duality

$$\mathcal{W}(\mathbb{P}, \mathbb{Q}) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} [f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}} [f(Y)]$$

$f: \mathcal{X} \rightarrow \mathbb{R}$ is a 1-Lipschitz *critic function*

$$\|f\|_{\text{Lip}} = \sup_{x,y \in \mathcal{X}} \frac{|f(x) - f(y)|}{\|x - y\|} = \sup_{x \in \mathcal{X}} \|\nabla f(x)\|$$



\mathcal{W} is continuous (bounded domain): if $\mathbb{Q}_\theta \rightarrow \mathbb{P}$, then $\mathcal{W}(\mathbb{Q}_\theta, \mathbb{P}) \rightarrow 0$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Idea: replace discriminator D_ψ with a critic f_ψ
- Need to enforce $\|f_\psi\|_{\text{Lip}} \leq 1$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Idea: replace discriminator D_ψ with a critic f_ψ
- Need to enforce $\|f_\psi\|_{\text{Lip}} \leq 1$...or, really, any constant

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Idea: replace discriminator D_ψ with a critic f_ψ
- Need to enforce $\|f_\psi\|_{\text{Lip}} \leq 1$...or, really, any constant
- $$f_\psi(x) = h_L(W_L h_{L-1}(\cdots h_1(W_1 x)))$$

so for usual deep nets,

$$\|f_\psi\|_{\text{Lip}} \leq \|h_L\|_{\text{Lip}} \|W_L\|_{\text{op}} \cdots \|h_1\|_{\text{Lip}} \|W_1\|_{\text{op}}$$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Idea: replace discriminator D_ψ with a critic f_ψ
- Need to enforce $\|f_\psi\|_{\text{Lip}} \leq 1$...or, really, any constant
- $$f_\psi(x) = h_L(W_L h_{L-1}(\cdots h_1(W_1 x)))$$

so for usual deep nets,

$$\|f_\psi\|_{\text{Lip}} \leq \|h_L\|_{\text{Lip}} \|W_L\|_{\text{op}} \cdots \|h_1\|_{\text{Lip}} \|W_1\|_{\text{op}}$$

- $\|W\|_{\text{op}} := \sup_{x \neq 0} \frac{\|Wx\|_2}{\|x\|} = \sigma_{\max}(W)$ is *spectral* norm

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Idea: replace discriminator D_ψ with a critic f_ψ
- Need to enforce $\|f_\psi\|_{\text{Lip}} \leq 1$...or, really, any constant
- $$f_\psi(x) = h_L(W_L h_{L-1}(\cdots h_1(W_1 x)))$$

so for usual deep nets,

$$\|f_\psi\|_{\text{Lip}} \leq \|h_L\|_{\text{Lip}} \|W_L\|_{\text{op}} \cdots \|h_1\|_{\text{Lip}} \|W_1\|_{\text{op}}$$

- $\|W\|_{\text{op}} := \sup_{x \neq 0} \frac{\|Wx\|_2}{\|x\|} = \sigma_{\max}(W)$ is *spectral* norm
- h is fixed; often $\|h\|_{\text{Lip}} \leq 1$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Need

$$\|f_{\psi}\|_{\text{Lip}} \leq \left(\prod_{j=1}^L \|h_j\|_{\text{Lip}} \right) \|W_L\|_{\text{op}} \cdots \|W_1\|_{\text{op}} \leq \text{const}$$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Need

$$\|f_{\psi}\|_{\text{Lip}} \leq \left(\prod_{j=1}^L \|h_j\|_{\text{Lip}} \right) \|W_L\|_{\text{op}} \cdots \|W_1\|_{\text{op}} \leq \text{const}$$

- WGANs: Bound $\|W_i\|_\infty \leq C$
 - This implies $\|W_i\|_{\text{op}} \leq \sqrt{k_i k_{i-1}} C$

WGAN [Arjovsky/Chintala/Bottou ICML-17]

- Need

$$\|f_{\psi}\|_{\text{Lip}} \leq \left(\prod_{j=1}^L \|h_j\|_{\text{Lip}} \right) \|W_L\|_{\text{op}} \cdots \|W_1\|_{\text{op}} \leq \text{const}$$

- WGANs: Bound $\|W_i\|_\infty \leq C$
 - This implies $\|W_i\|_{\text{op}} \leq \sqrt{k_i k_{i-1}} C$
- ...this turns out not to be a great idea.

WGAN-GP [Gulrajani+ NeurIPS-17]

- Controlling $\|\nabla f(X)\|$ *everywhere* is hard

WGAN-GP [Gulrajani+ NeurIPS-17]

- Controlling $\|\nabla f(X)\|$ *everywhere* is hard
- Instead, control $\|\nabla f(\tilde{X})\|$ *on average, near the data*

WGAN-GP [Gulrajani+ NeurIPS-17]

- Controlling $\|\nabla f(X)\|$ *everywhere* is hard
- Instead, control $\|\nabla f(\tilde{X})\|$ *on average, near the data*

$$\mathbb{E}_{\tilde{X} \sim \mathbb{S}} (\|\nabla_{\tilde{X}} f_{\psi}(\tilde{X})\| - 1)^2, \quad \mathbb{S} \text{ “between” } \mathbb{P} \text{ and } Q_{\theta}$$

WGAN-GP [Gulrajani+ NeurIPS-17]

- Controlling $\|\nabla f(\mathbf{X})\|$ *everywhere* is hard
- Instead, control $\|\nabla f(\tilde{\mathbf{X}})\|$ *on average, near the data*

$$\mathbb{E}_{\tilde{\mathbf{X}} \sim \mathbb{S}} \left(\|\nabla_{\tilde{\mathbf{X}}} f_{\psi}(\tilde{\mathbf{X}})\| - 1 \right)^2, \quad \mathbb{S} \text{ “between” } \mathbb{P} \text{ and } \mathbb{Q}_{\theta}$$

- Specifically: $\tilde{\mathbf{X}} = \theta \mathbf{X} + (1 - \theta) \mathbf{Y}, \theta \sim \text{Uniform}([0, 1])$

WGAN-GP [Gulrajani+ NeurIPS-17]

- Controlling $\|\nabla f(\mathbf{X})\|$ *everywhere* is hard
- Instead, control $\|\nabla f(\tilde{\mathbf{X}})\|$ *on average, near the data*

$$\mathbb{E}_{\tilde{\mathbf{X}} \sim \mathbb{S}} \left(\|\nabla_{\tilde{\mathbf{X}}} f_{\psi}(\tilde{\mathbf{X}})\| - 1 \right)^2, \quad \mathbb{S} \text{ “between” } \mathbb{P} \text{ and } \mathbb{Q}_{\theta}$$

- Specifically: $\tilde{\mathbf{X}} = \theta \mathbf{X} + (1 - \theta) \mathbf{Y}, \theta \sim \text{Uniform}([0, 1])$
- Works well!

WGAN-GP [Gulrajani+ NeurIPS-17]

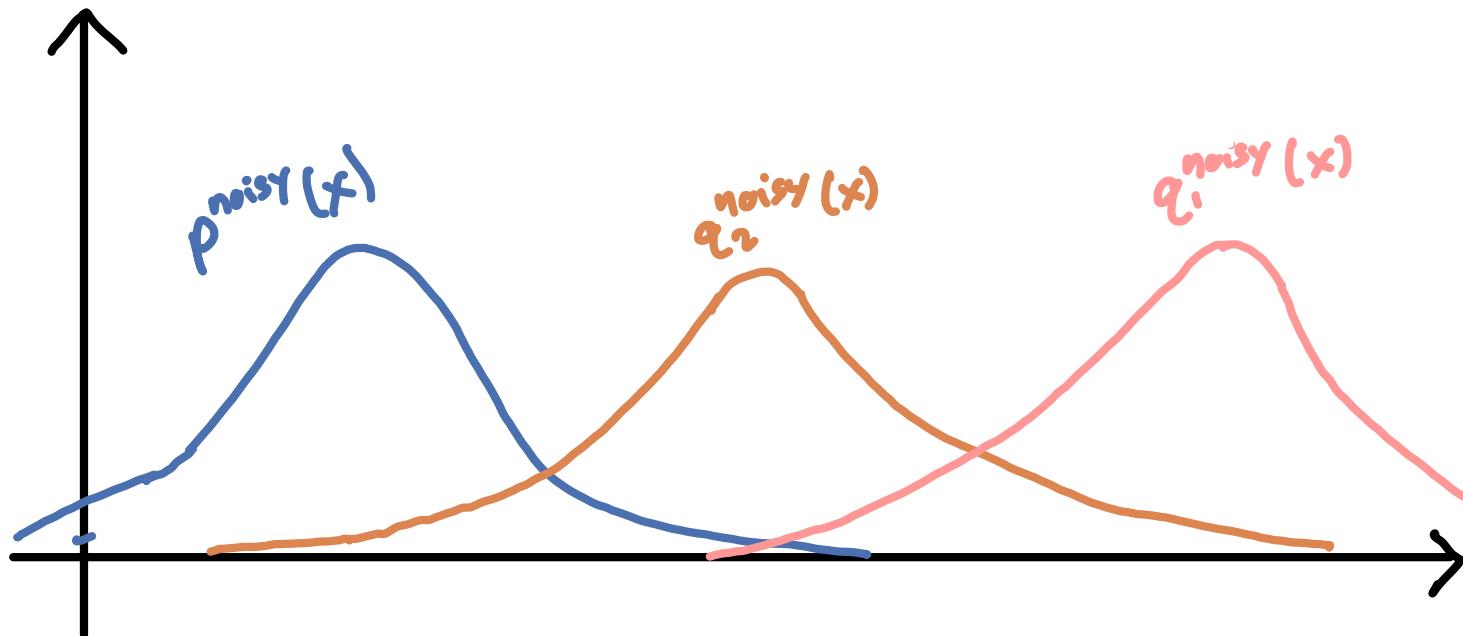
- Controlling $\|\nabla f(\mathbf{X})\|$ *everywhere* is hard
- Instead, control $\|\nabla f(\tilde{\mathbf{X}})\|$ *on average, near the data*

$$\mathbb{E}_{\tilde{\mathbf{X}} \sim \mathbb{S}} (\|\nabla_{\tilde{\mathbf{X}}} f_{\psi}(\tilde{\mathbf{X}})\| - 1)^2, \quad \mathbb{S} \text{ “between” } \mathbb{P} \text{ and } \mathbb{Q}_{\theta}$$

- Specifically: $\tilde{\mathbf{X}} = \theta \mathbf{X} + (1 - \theta) \mathbf{Y}, \theta \sim \text{Uniform}([0, 1])$
- Works well! But...does it really estimate Wasserstein?

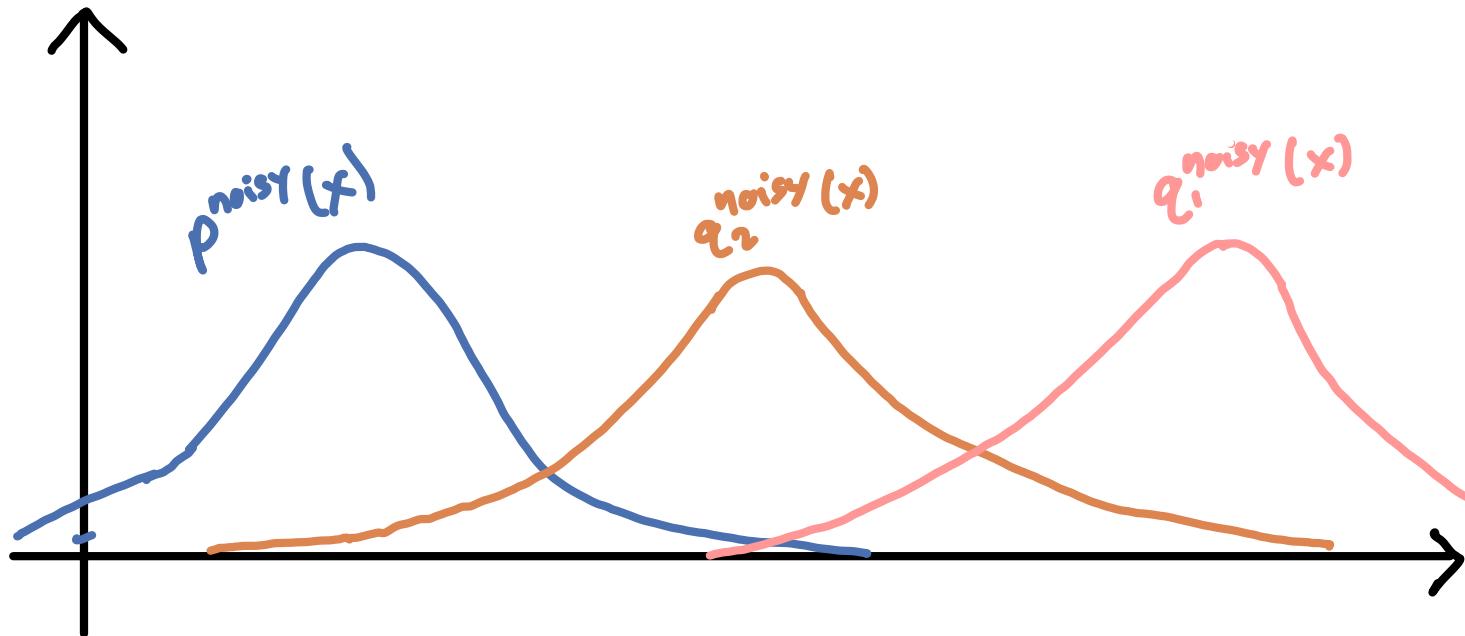
Solution 2: add noise

- Can keep JS if we “blur” the data
- Use $\mathbf{X} + \varepsilon, \mathbf{Y} + \varepsilon'$ for some independent, full-dim noise ε



Solution 2: add noise

- Can keep JS if we “blur” the data
- Use $\mathbf{X} + \varepsilon, \mathbf{Y} + \varepsilon'$ for some independent, full-dim noise ε



- But...how much noise ε to add? Also need more samples.

Solution 2: add noise in closed form

- Use $\mathbf{X} + \varepsilon, \mathbf{Y} + \varepsilon'$ for some independent, full-dim noise ε
- If $\varepsilon \sim \mathcal{N}(0, \gamma I)$ and $\gamma \rightarrow 0$, get [Mescheder+ NeurIPS-17]

$$\gamma \mathbb{E}_{\mathbb{P}}[(1 - D_{\psi})^2 \|\nabla \log(D_{\psi})\|^2] + \gamma \mathbb{E}_{Q_{\theta}}[D_{\psi}^2 \|\nabla \log(D_{\psi})\|^2]$$

Solution 2: add noise in closed form

- Use $\mathbf{X} + \varepsilon, \mathbf{Y} + \varepsilon'$ for some independent, full-dim noise ε
- If $\varepsilon \sim \mathcal{N}(0, \gamma I)$ and $\gamma \rightarrow 0$, get [Mescheder+ NeurIPS-17]

$$\gamma \mathbb{E}_{\mathbb{P}}[(1 - D_{\psi})^2 \|\nabla \log(D_{\psi})\|^2] + \gamma \mathbb{E}_{Q_{\theta}}[D_{\psi}^2 \|\nabla \log(D_{\psi})\|^2]$$

- Similar kind of gradient penalty!

Solution 2: add noise in closed form

- Use $\mathbf{X} + \varepsilon, \mathbf{Y} + \varepsilon'$ for some independent, full-dim noise ε
- If $\varepsilon \sim \mathcal{N}(0, \gamma I)$ and $\gamma \rightarrow 0$, get [Mescheder+ NeurIPS-17]

$$\gamma \mathbb{E}_{\mathbb{P}}[(1 - D_{\psi})^2 \|\nabla \log(D_{\psi})\|^2] + \gamma \mathbb{E}_{Q_{\theta}}[D_{\psi}^2 \|\nabla \log(D_{\psi})\|^2]$$

- Similar kind of gradient penalty!
- Can also simplify to e.g. [Mescheder+ ICML-18]

$$\gamma \mathbb{E}_{\mathbf{X} \sim \mathbb{P}}[\|\nabla D_{\psi}(\mathbf{X})\|^2]$$

Solution 3: Spectral norm [Miyato+ ICLR-18]

- Regular deep nets: $f_\ell = h_\ell(W_\ell f_{\ell-1}(x))$
- Spectral normalization: $f_\ell = h\left(\frac{1}{\|W_\ell\|_{\text{op}}} W_\ell f_{\ell-1}(x)\right)$
- Guarantees $\|f\|_{\text{Lip}} \leq 1$, like WGAN
- With tricks, faster to evaluate than gradient penalties
- Not as well understood yet

Solution 3: Spectral norm [Miyato+ ICLR-18]

- Regular deep nets: $f_\ell = h_\ell(W_\ell f_{\ell-1}(x))$
- Spectral normalization: $f_\ell = h\left(\frac{1}{\|W_\ell\|_{\text{op}}} W_\ell f_{\ell-1}(x)\right)$
- Guarantees* $\|f\|_{\text{Lip}} \leq 1$, like WGAN
- With tricks, faster to evaluate than gradient penalties
- Not as well understood yet

New samples [Mescheder+ ICML-18]



How to evaluate?



FID [Heusel+ NeurIPS-17] and KID [Bińkowski+ ICLR-18]

- Consider distance between distributions of image features
- Features $\phi(x)$ from a pretrained ImageNet classifier

FID [Heusel+ NeurIPS-17] and KID [Bińkowski+ ICLR-18]

- Consider distance between distributions of image features
- Features $\phi(x)$ from a pretrained ImageNet classifier
- FID: $\|\mu_{\mathbb{P}} - \mu_{Q_\theta}\|^2 + \text{Tr} \left(\Sigma_{\mathbb{P}} + \Sigma_{Q_\theta} - 2(\Sigma_{\mathbb{P}} \Sigma_{Q_\theta})^{\frac{1}{2}} \right)$

FID [Heusel+ NeurIPS-17] and KID [Bińkowski+ ICLR-18]

- Consider distance between distributions of image features
- Features $\phi(x)$ from a pretrained ImageNet classifier
- FID: $\|\mu_{\mathbb{P}} - \mu_{Q_\theta}\|^2 + \text{Tr} \left(\Sigma_{\mathbb{P}} + \Sigma_{Q_\theta} - 2(\Sigma_{\mathbb{P}} \Sigma_{Q_\theta})^{\frac{1}{2}} \right)$
 - Estimator very biased, small variance

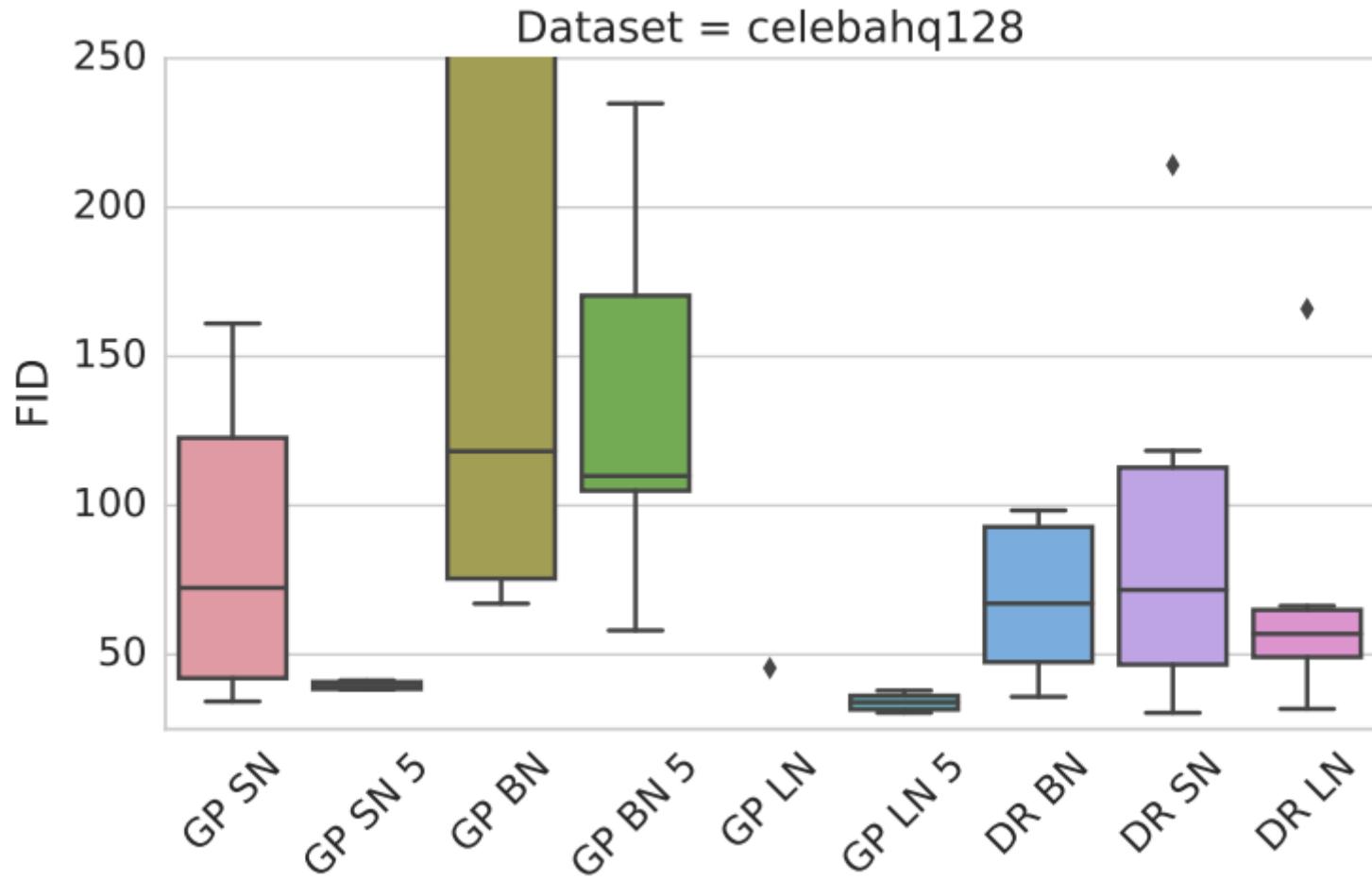
FID [Heusel+ NeurIPS-17] and KID [Bińkowski+ ICLR-18]

- Consider distance between distributions of image features
- Features $\phi(x)$ from a pretrained ImageNet classifier
- FID: $\|\mu_{\mathbb{P}} - \mu_{Q_\theta}\|^2 + \text{Tr} \left(\Sigma_{\mathbb{P}} + \Sigma_{Q_\theta} - 2(\Sigma_{\mathbb{P}} \Sigma_{Q_\theta})^{\frac{1}{2}} \right)$
 - Estimator very biased, small variance
- KID: use Maximum Mean Discrepancy instead

FID [Heusel+ NeurIPS-17] and KID [Bińkowski+ ICLR-18]

- Consider distance between distributions of image features
- Features $\phi(x)$ from a pretrained ImageNet classifier
- FID: $\|\mu_{\mathbb{P}} - \mu_{Q_\theta}\|^2 + \text{Tr} \left(\Sigma_{\mathbb{P}} + \Sigma_{Q_\theta} - 2(\Sigma_{\mathbb{P}} \Sigma_{Q_\theta})^{\frac{1}{2}} \right)$
 - Estimator very biased, small variance
- KID: use Maximum Mean Discrepancy instead
 - Similar distance with unbiased, \sim normal estimator!

Comparing approaches [Kurach+ ICML-19]



Maximum Mean Discrepancy

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

$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

Maximum Mean Discrepancy

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

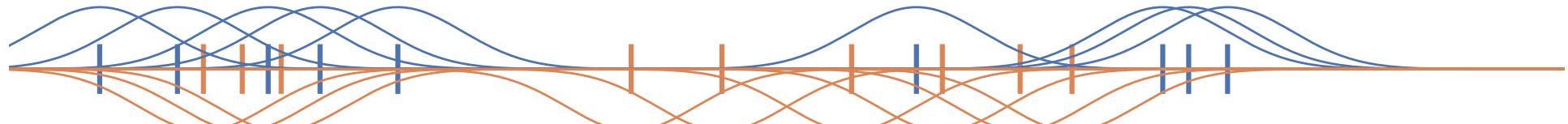
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

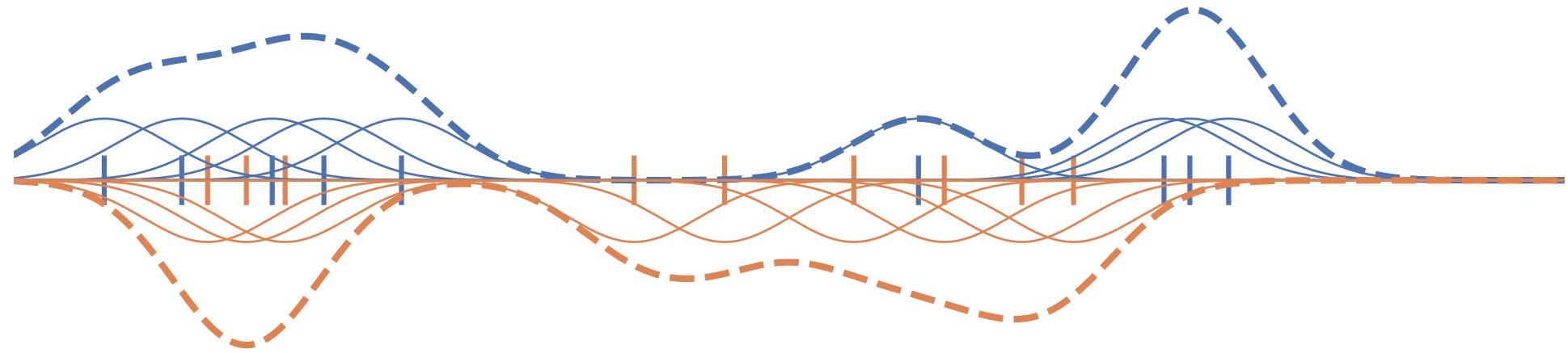
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

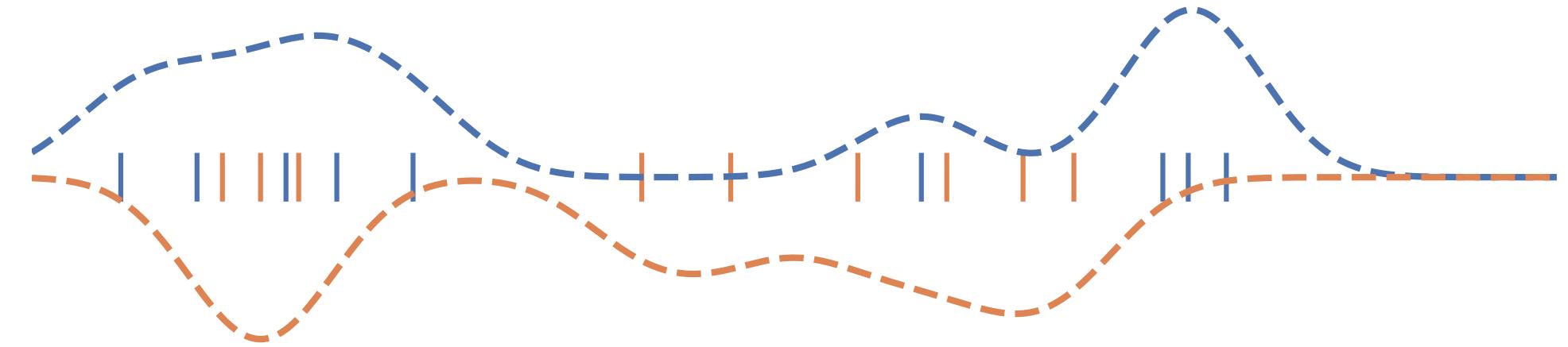
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

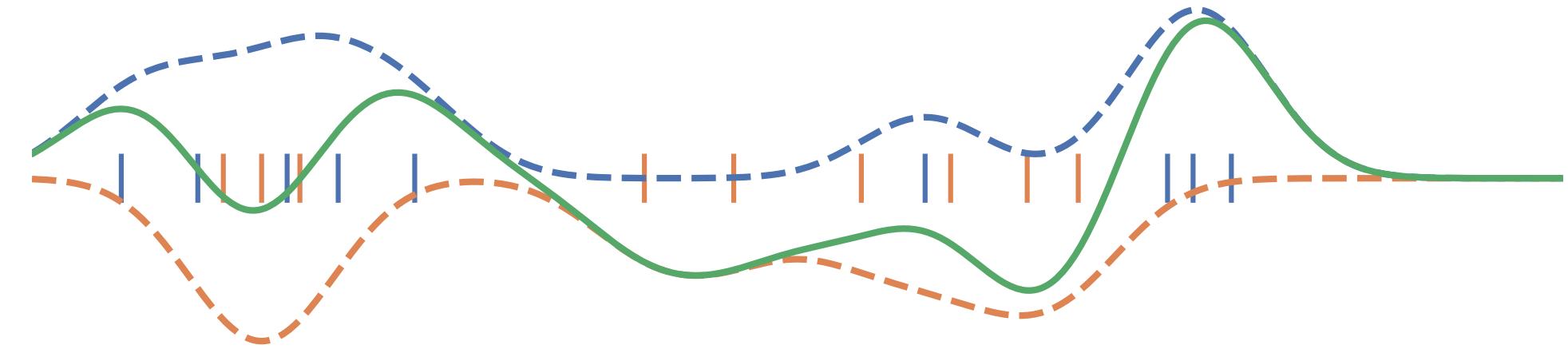
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

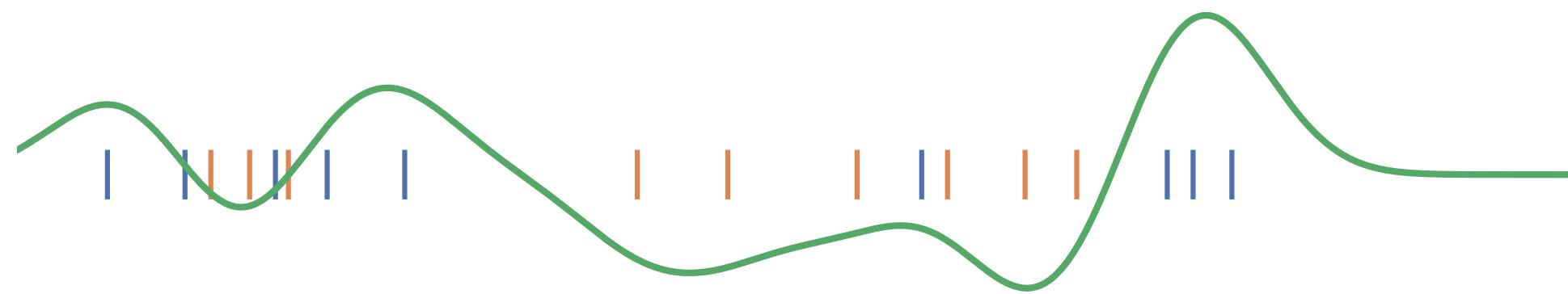
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

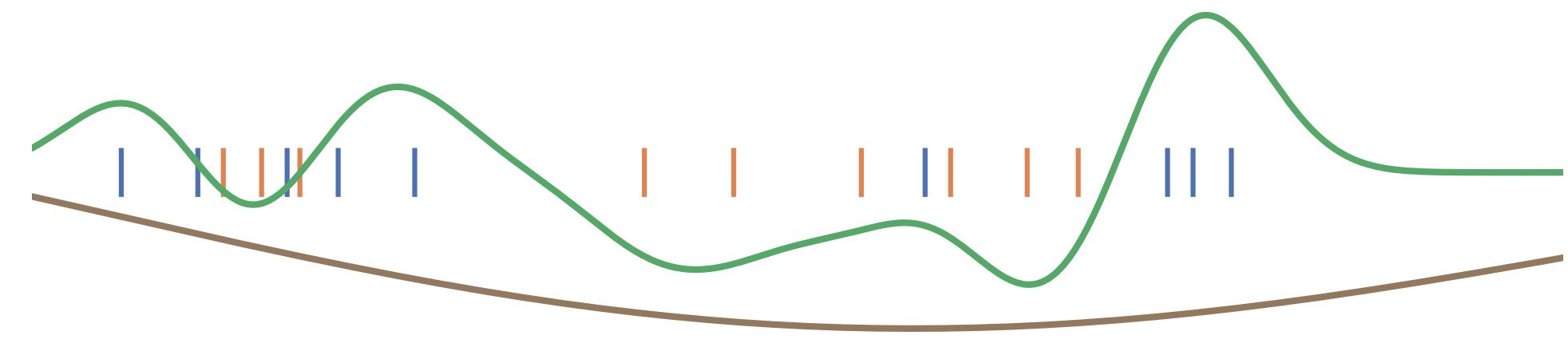
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

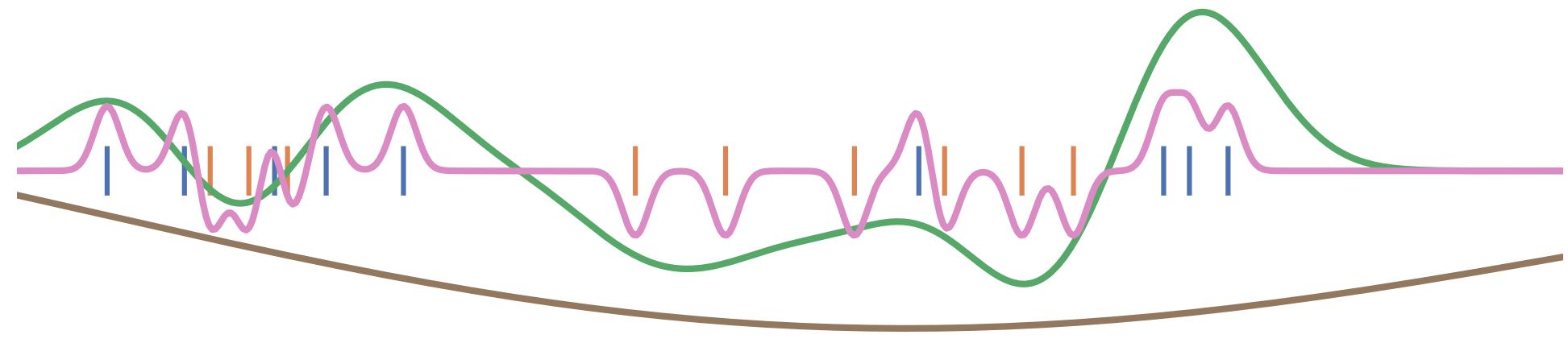
$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Maximum Mean Discrepancy

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

$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

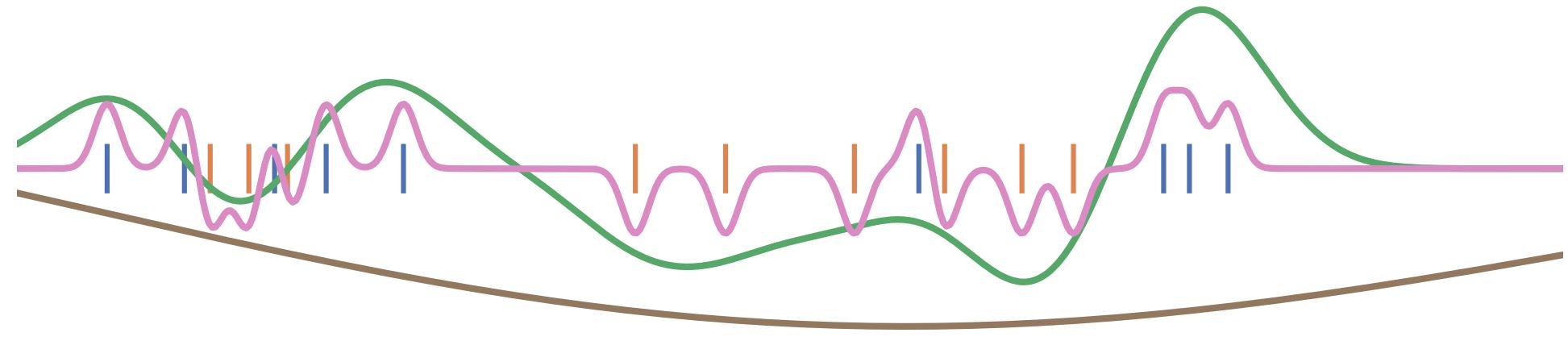


Maximum Mean Discrepancy

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

$\|f\|_{\mathcal{H}_k}$ is smoothness induced by kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

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



Estimating MMD

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

Estimating MMD

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

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

Estimating MMD

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

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

K_{XX}

—	1.0	0.2	0.6
—	0.2	1.0	0.5
—	0.6	0.5	1.0

Estimating MMD

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

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

K_{XX}

K_{YY}

1.0	0.2	0.6
0.2	1.0	0.5
0.6	0.5	1.0

1.0	0.8	0.7
0.8	1.0	0.6
0.7	0.6	1.0

Estimating MMD

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

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

K_{XX}

1.0	0.2	0.6
0.2	1.0	0.5
0.6	0.5	1.0

K_{YY}

1.0	0.8	0.7
0.8	1.0	0.6
0.7	0.6	1.0

K_{XY}

0.3	0.1	0.2
0.2	0.3	0.3
0.2	0.1	0.4

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

Generator (Q_θ)



Critic



MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

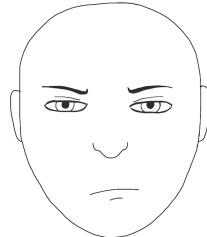
Generator (Q_θ)



Critic



How are these?



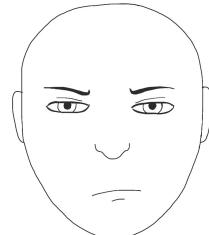
MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

Generator (Q_θ)



How are these?



Critic



Target (\mathbb{P})



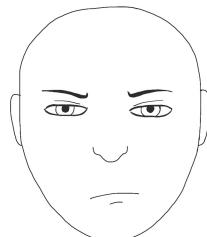
MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

Generator (Q_θ)



How are these?



Critic



Target (\mathbb{P})

Not great! $\widehat{\text{MMD}}(Q_\theta, \mathbb{P}) = 0.75$

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

Generator (Q_θ)



How are these?



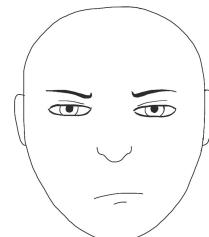
Critic



Target (\mathbb{P})



Not great! $\widehat{\text{MMD}}(Q_\theta, \mathbb{P}) = 0.75$



:(| I'll try harder...

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

- No need for a discriminator – just minimize $\widehat{\text{MMD}}_k$!
- Continuous loss

Generator (Q_θ)



How are these?



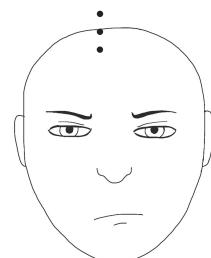
Critic



Target (\mathbb{P})



Not great! $\widehat{\text{MMD}}(Q_\theta, \mathbb{P}) = 0.75$



:(| I'll try harder...

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

MNIST, mix of Gaussian kernels



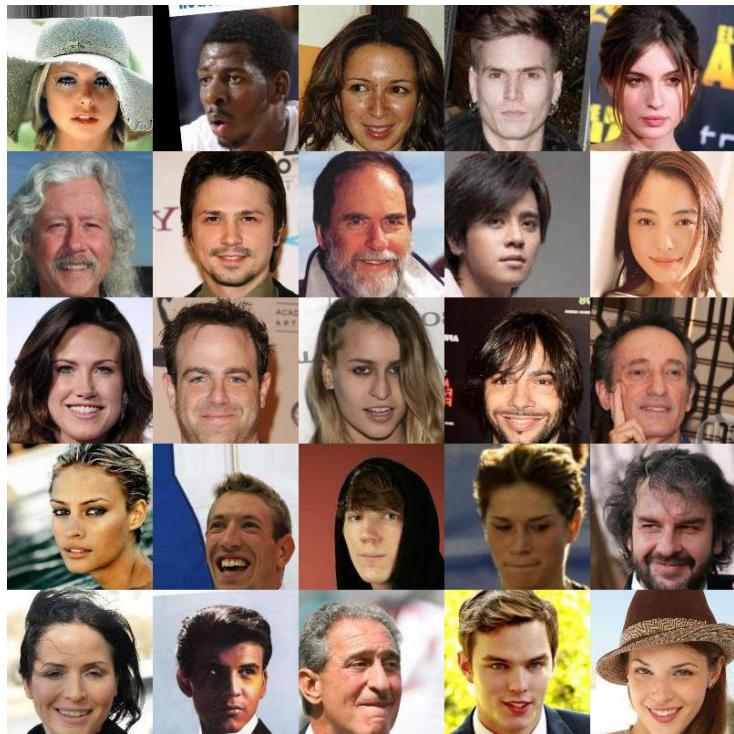
P

MMD models [Li+ ICML-15, Dziugaite+ UAI-15]

MNIST, mix of Gaussian kernels

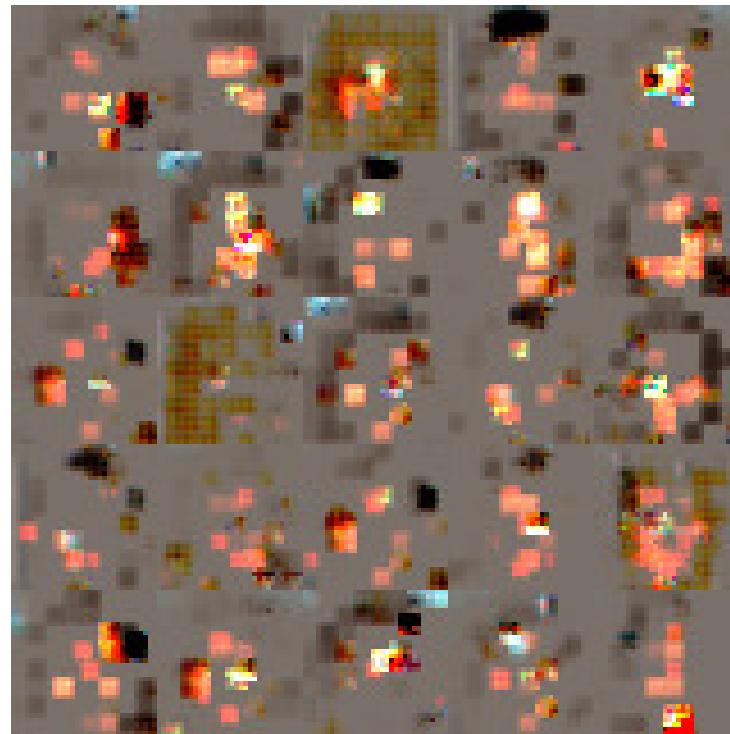
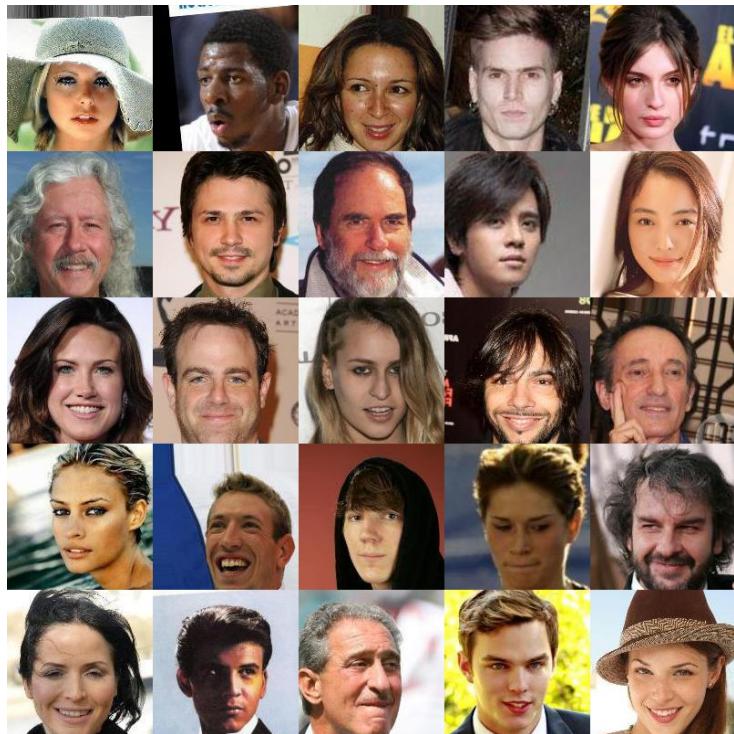


Celeb-A, mix of rational quadratic + linear kernels



P

Celeb-A, mix of rational quadratic + linear kernels



\mathbb{P}

Q_θ

MMD loss with a smarter kernel

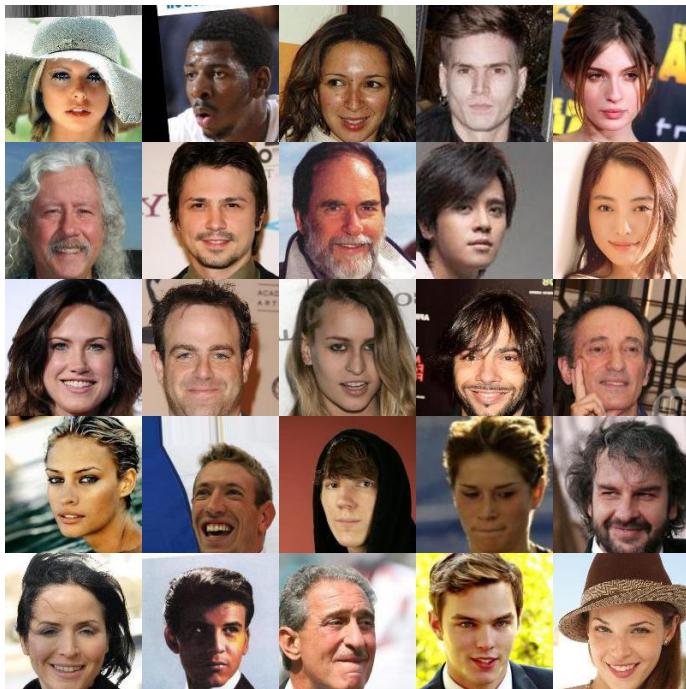
$$k(x, y) = k_{\text{top}}(\phi(x), \phi(y))$$

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

MMD loss with a smarter kernel

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

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

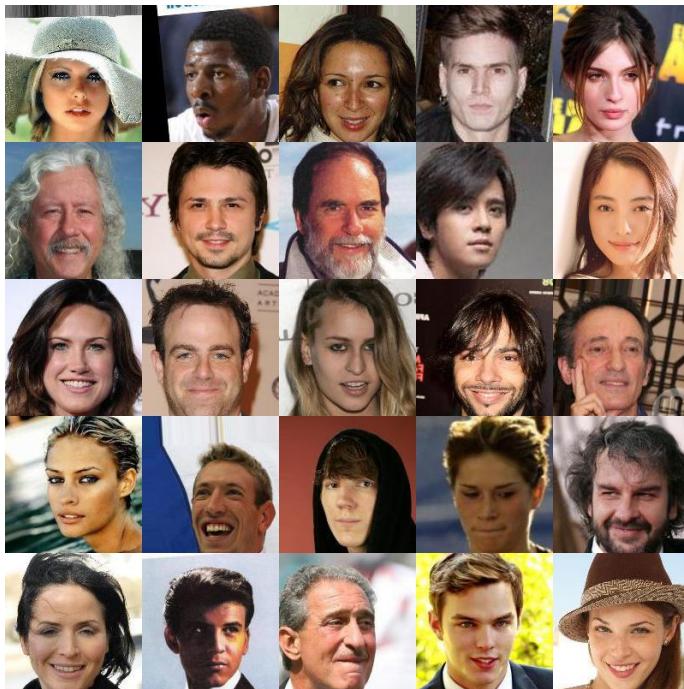


P

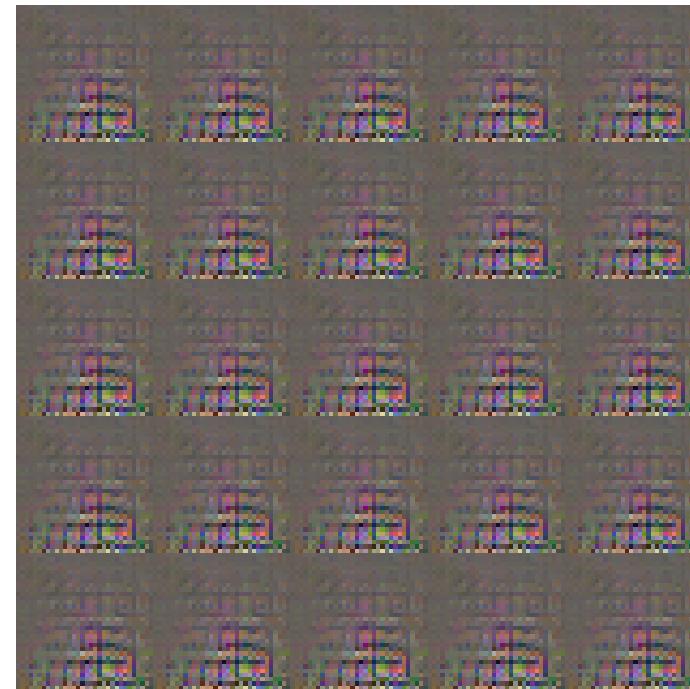
MMD loss with a smarter kernel

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

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



\mathbb{P}



Q_θ

MMD loss with a smarter kernel

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

- ϕ : .
- k_{top}

We just got adversarial examples!



88% **tabby cat**

adversarial
perturbation



99% **guacamole**

polynomial



[anishathalye/obfuscated-gradients]



\mathbb{P}



Q_θ

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

- Don't just use one kernel, use a *class* parameterized by ψ :

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

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

- Don't just use one kernel, use a *class* parameterized by ψ :

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

- New distance based on *all* these kernels:

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

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

- Don't just use one kernel, use a *class* parameterized by ψ :

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

- New distance based on *all* these kernels:

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

- Turns out that \mathcal{D}_{MMD} *isn't* continuous: have $\mathbb{Q}_{\theta} \rightarrow \mathbb{P}$ but $\mathcal{D}_{\text{MMD}}(\mathbb{Q}_{\theta}, \mathbb{P}) \not\rightarrow 0$

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

- Don't just use one kernel, use a *class* parameterized by ψ :

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

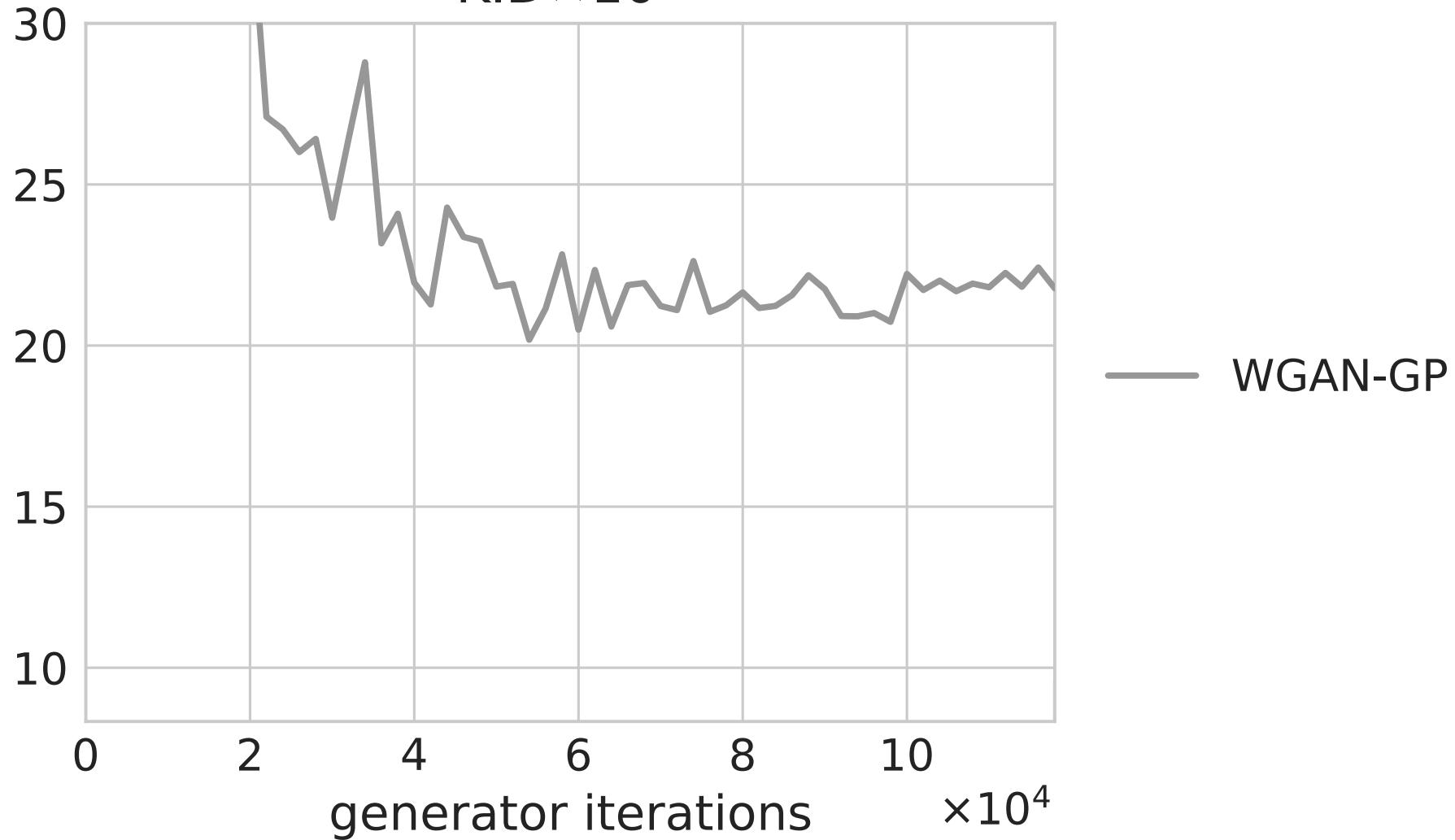
- New distance based on *all* these kernels:

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

- Turns out that \mathcal{D}_{MMD} *isn't* continuous: have $\mathbb{Q}_{\theta} \rightarrow \mathbb{P}$ but $\mathcal{D}_{\text{MMD}}(\mathbb{Q}_{\theta}, \mathbb{P}) \not\rightarrow 0$
- Scaled MMD GANs [Arbel+ NeurIPS-18] correct \mathcal{D}_{MMD} with a gradient penalty to make it continuous

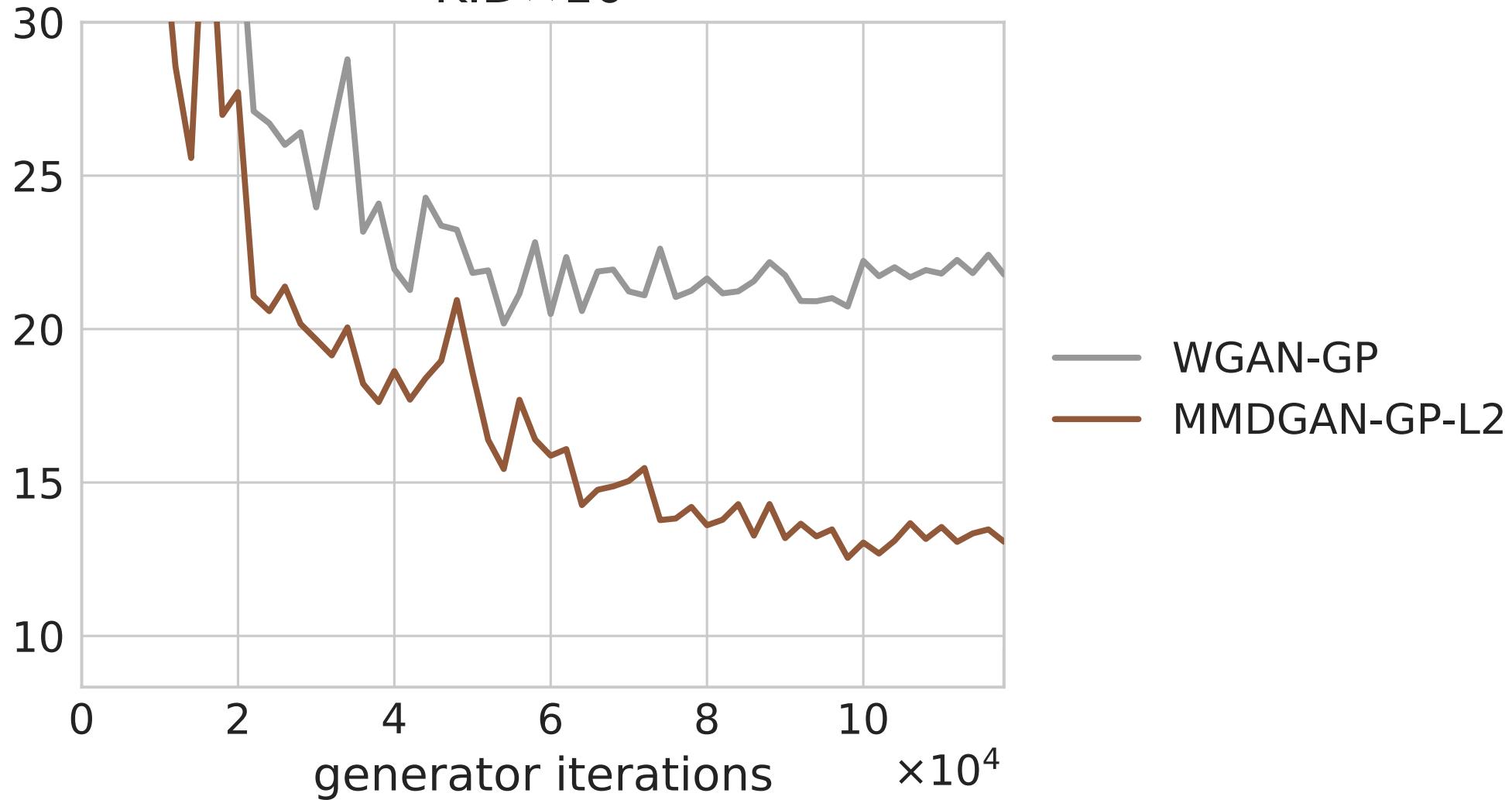
Why MMD GANs?

- “Easy parts” of the optimization done *in closed form*



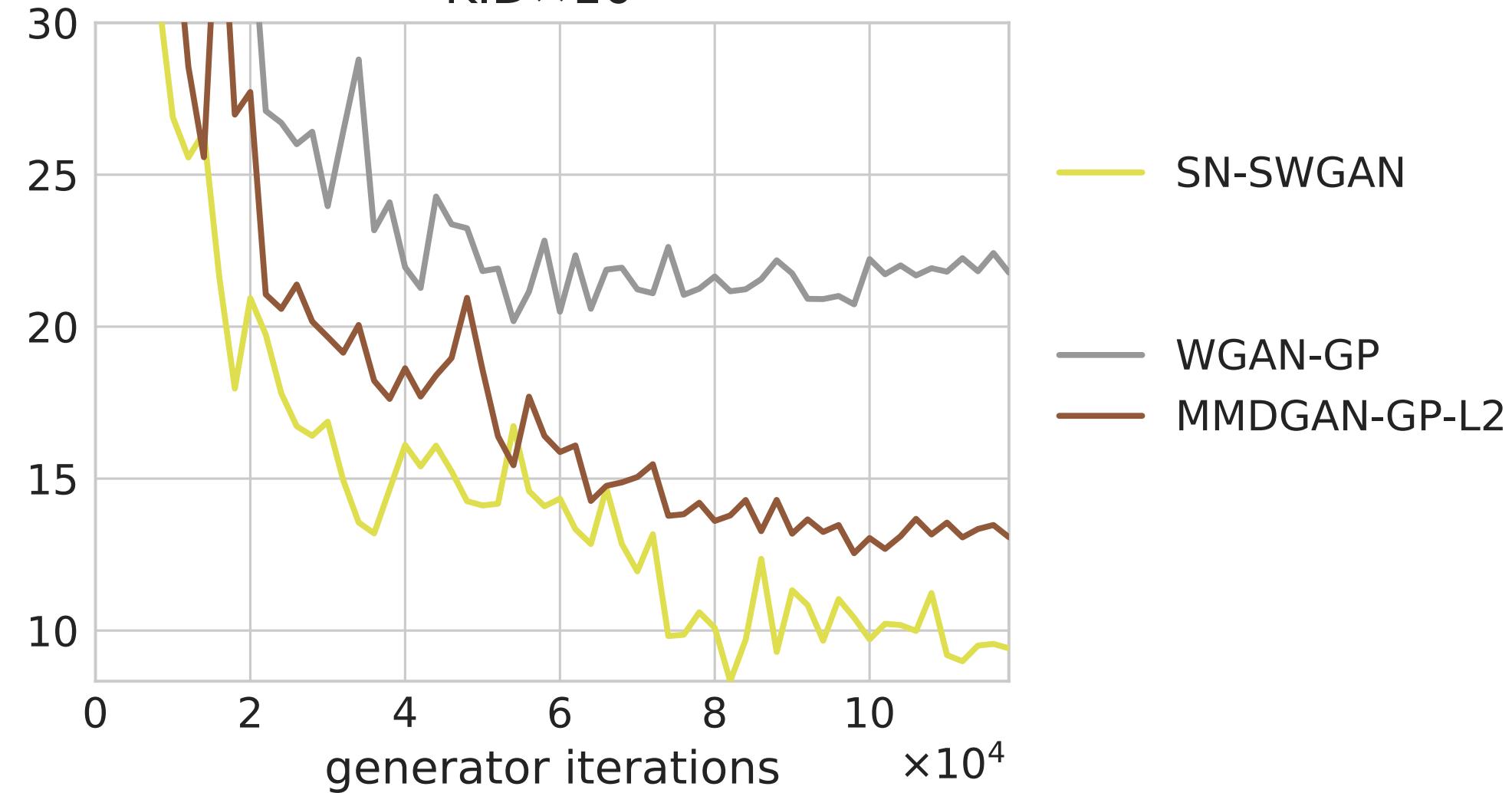
Why MMD GANs?

- “Easy parts” of the optimization done *in closed form*



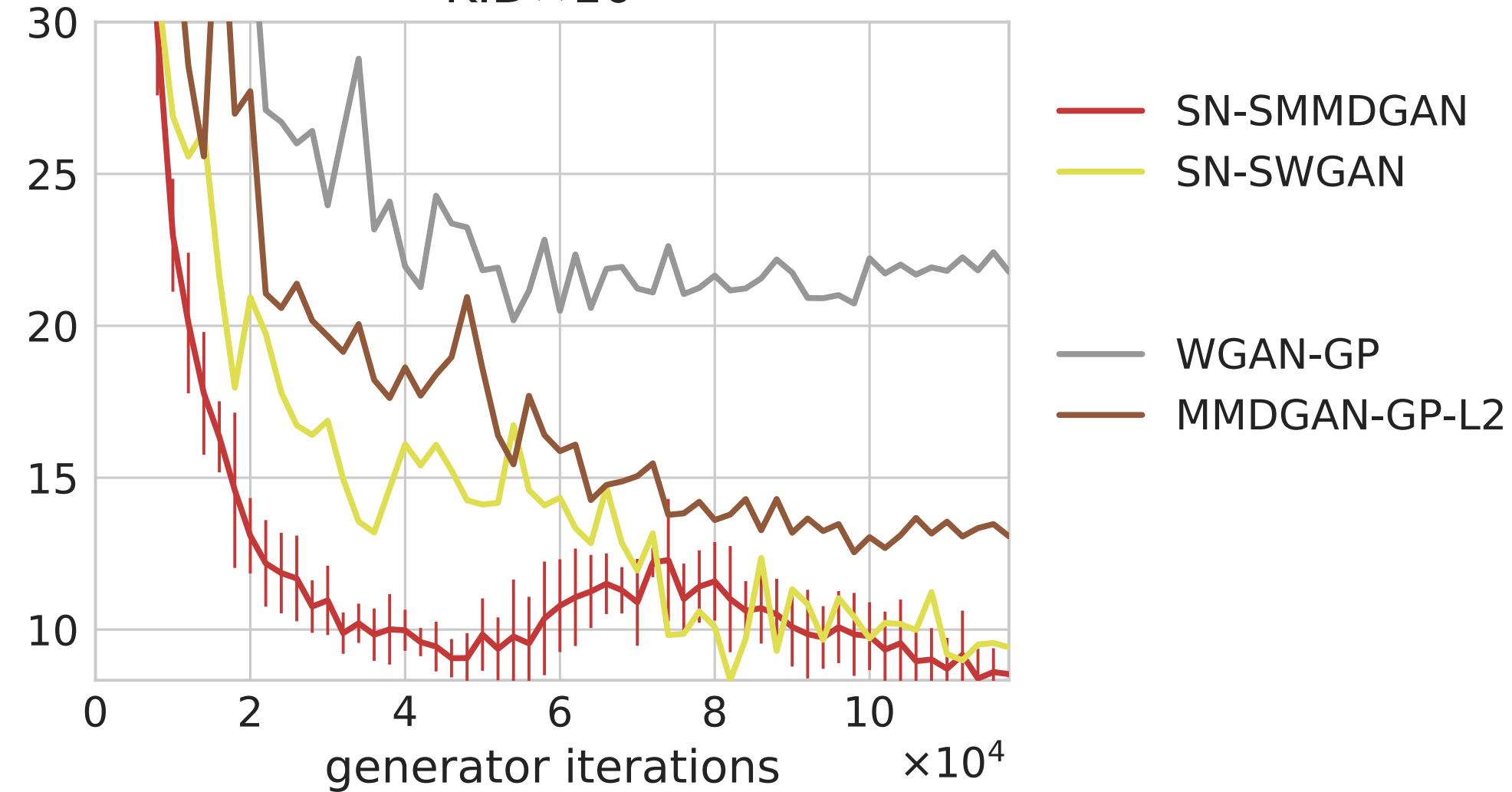
Why MMD GANs?

- “Easy parts” of the optimization done *in closed form*



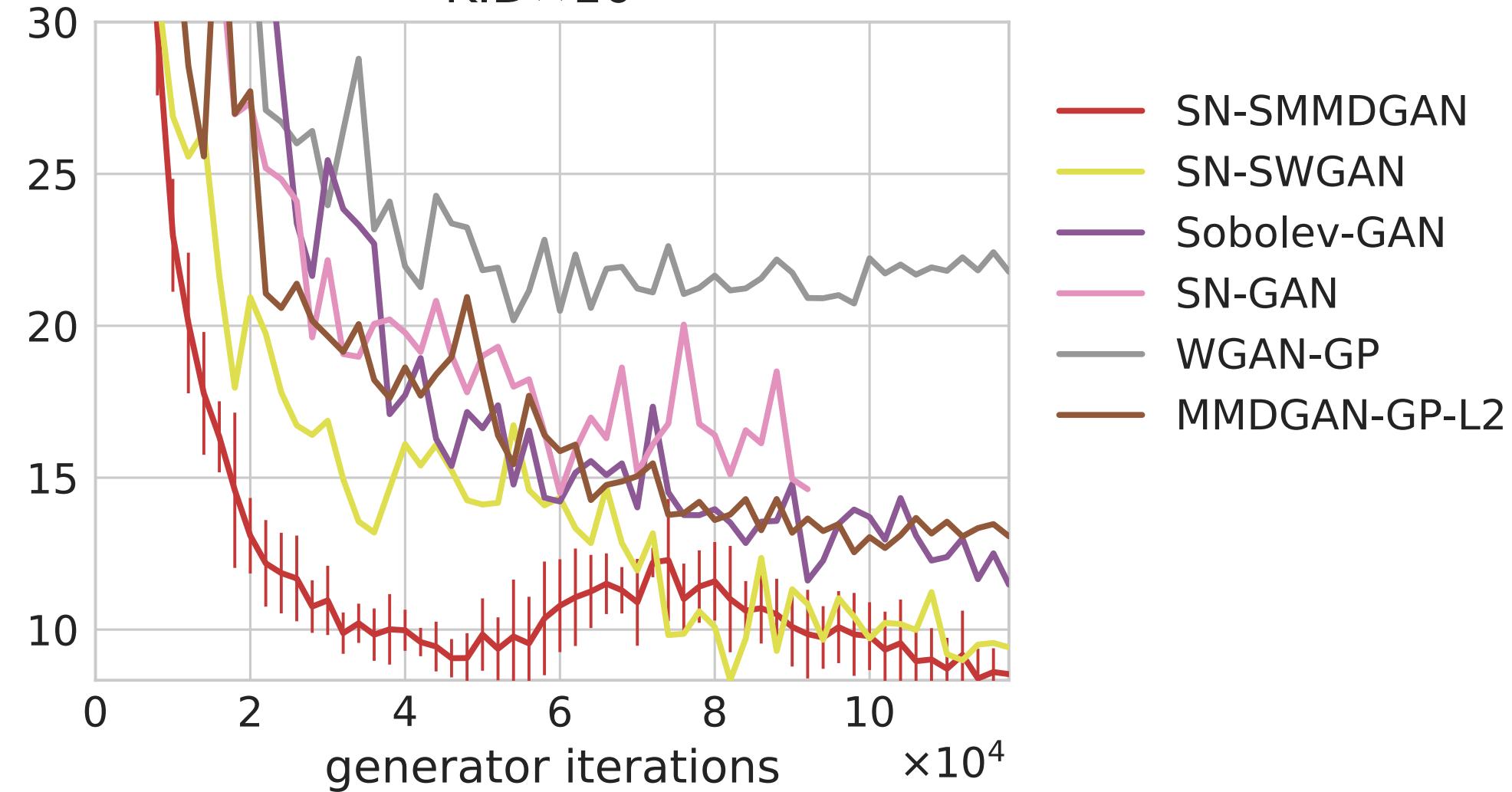
Why MMD GANs?

- “Easy parts” of the optimization done *in closed form*

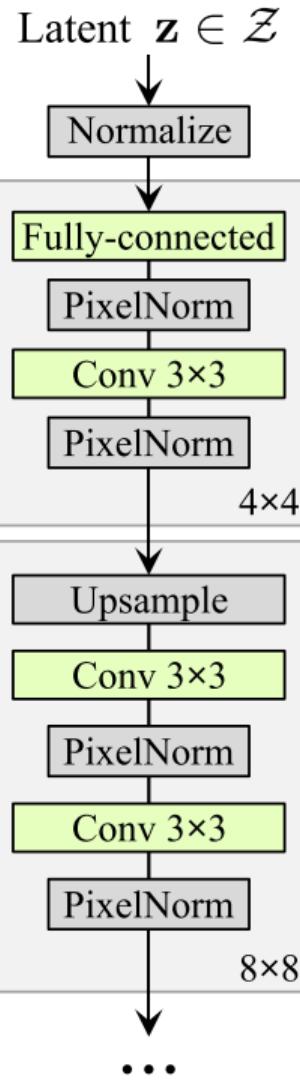


Why MMD GANs?

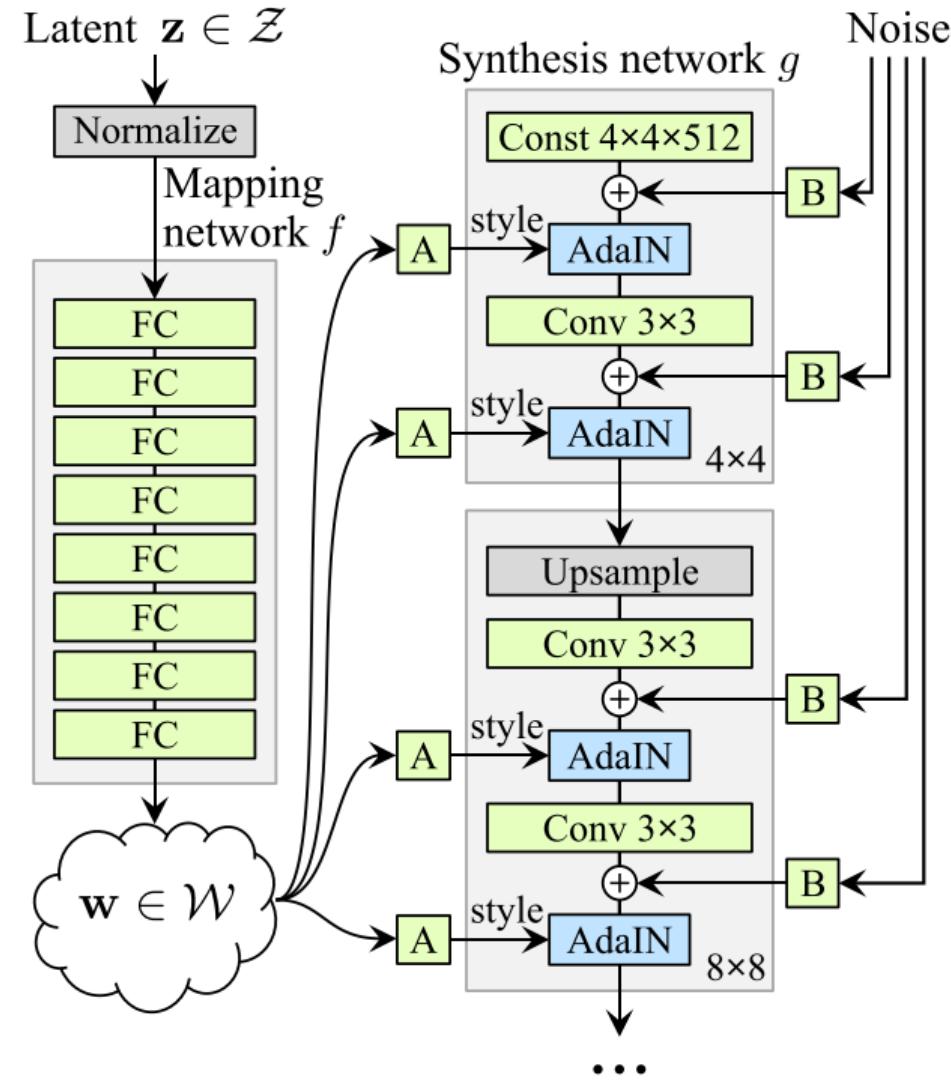
- “Easy parts” of the optimization done *in closed form*



StyleGANs [Karras+ 2018]



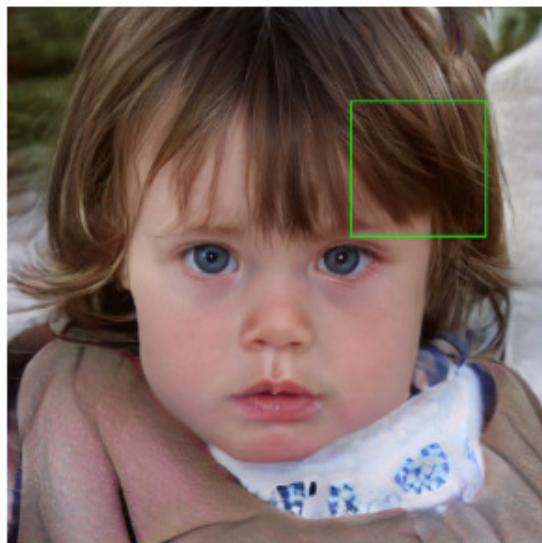
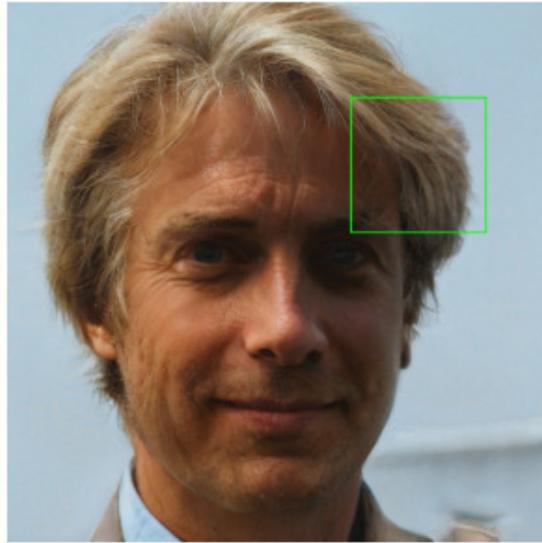
(a) Traditional



(b) Style-based generator

StyleGAN: latent structure

StyleGAN: local noise



(a) Generated image

(b) Stochastic variation

(c) Standard deviation

StyleGANs on a different domain [[@roadrunning01](#)]



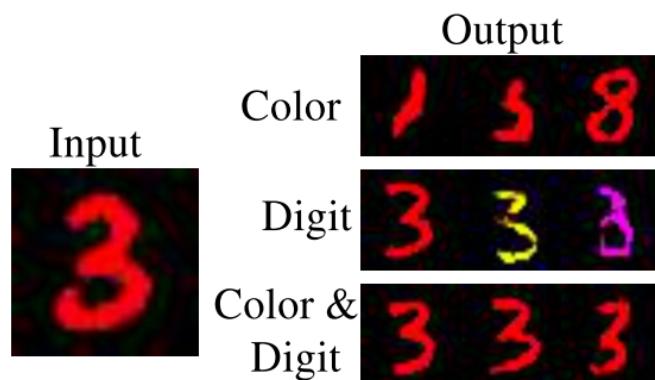
Finding samples you want [Jitkrittum+ ICML-19]

If we want to find “more samples like $\{\mathbf{X}\}$ ”:

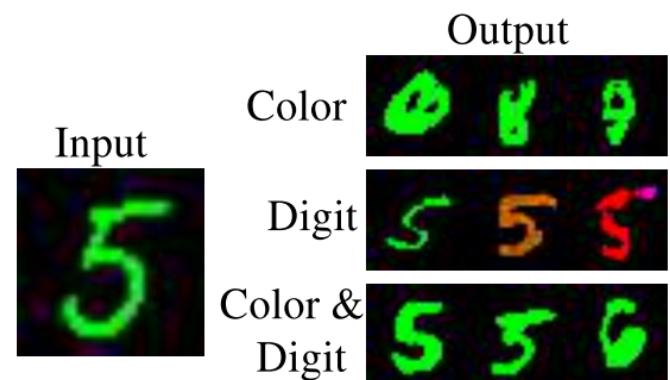
$$\min_{\{\mathbf{Z}_1, \dots, \mathbf{Z}_n\}} \widehat{\text{MMD}}_k^2(\{\mathbf{X}_i\}_{i=1}^m, \{G_\theta(\mathbf{Z}_i)\}_{i=1}^n)$$



(a) Samples from DCGAN

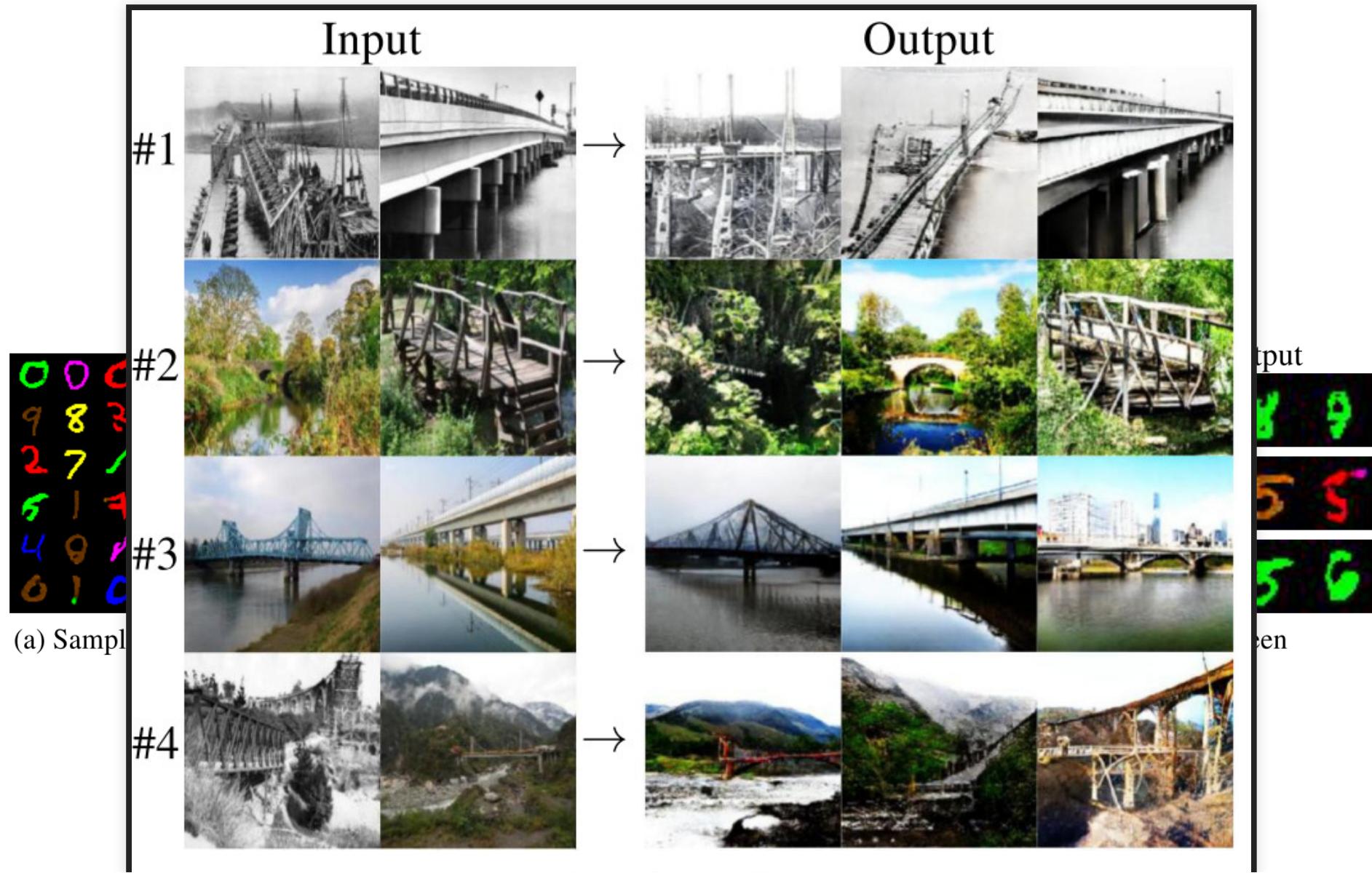


(b) Input: digit 3 in red



(c) Input: digit 5 in green

Finding samples you want [Jitkrittum+ ICML-19]



Conditional GANs and BigGAN

- Conditional GANs: [Mirza+ 2014]
 - Just add a class label as input to G_θ and D_ψ
- BigGAN [Brock+ ICLR-19]: a bunch of tricks to make it huge

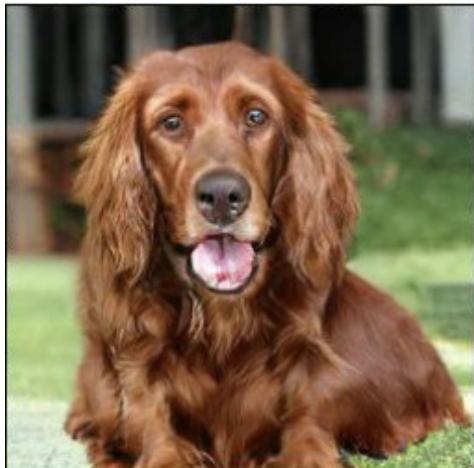


Image-to-image translation [Isola+ CVPR-17]

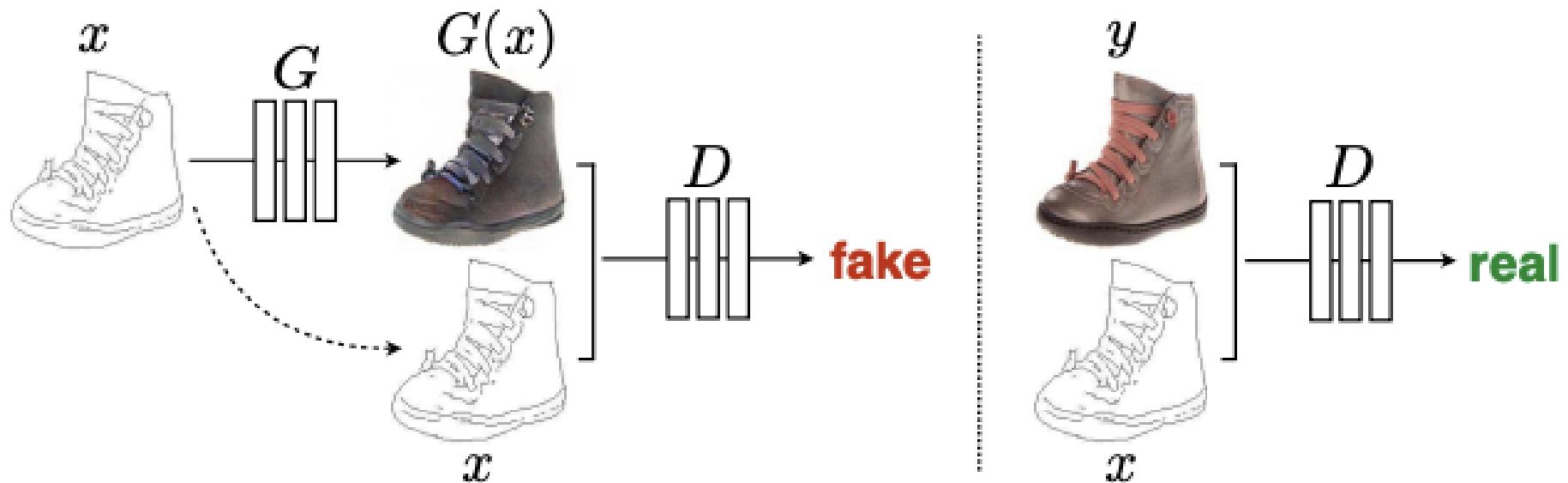
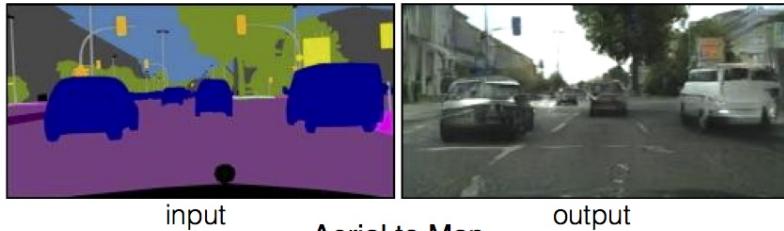


Figure 2: Training a conditional GAN to map edges→photo. The discriminator, D , learns to classify between fake (synthesized by the generator) and real {edge, photo} tuples. The generator, G , learns to fool the discriminator. Unlike an unconditional GAN, both the generator and discriminator observe the input edge map.

Image-to-image translation [Isola+ CVPR-17]

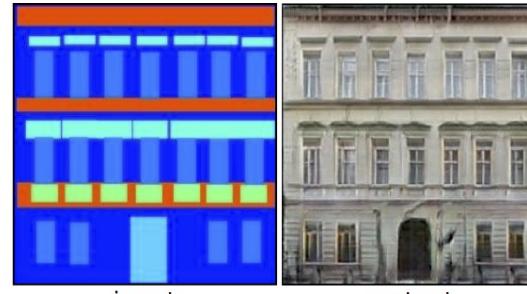
Labels to Street Scene



input

output

Labels to Facade



input

output

BW to Color



input

output

Aerial to Map



input

output

Day to Night



input

output

Edges to Photo



input

output

CycleGAN [Zhu+ ICCV-17]

Monet \curvearrowright Photos



Monet \rightarrow photo

Zebras \curvearrowright Horses



zebra \rightarrow horse

Summer \curvearrowright Winter



summer \rightarrow winter

photo \rightarrow Monet



horse \rightarrow zebra

winter \rightarrow summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e

Pose-to-image translation [Chan+ 2018]

▪

YouTube: Everybody Dance Now

DeepFakes

YouTube: Mark Zuckerberg 'deepfake' will remain online

More

- Optimal transport stuff:
 - Gabriel Peyré: *Optimal transport for machine learning* talk
 - Peyré and Cuturi, [Computational Optimal Transport](#) book
 - Kantorovich Initiative: kantorovich.org
 - [Pacific Interdisciplinary Hub on Optimal Transport](#)
- GANs / generative models...so much.