

Better gradient regularisation for MMD GANs

Arthur Gretton

Gatsby Computational Neuroscience Unit,
University College London

Theoretical Foundations and Applications of Deep Generative
Models, ICML 2018

Comparing two samples

- Given: Samples from unknown distributions P and Q .
- Goal: do P and Q differ?



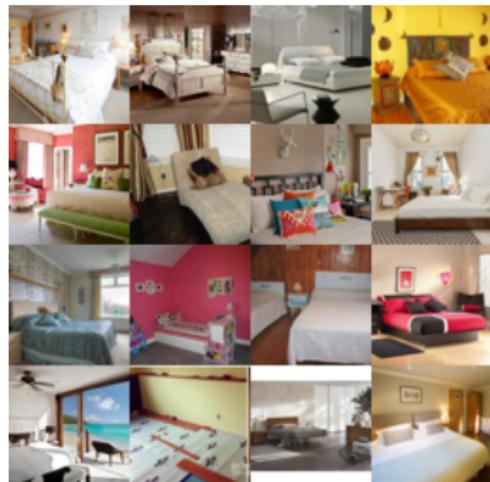
$\sim P$



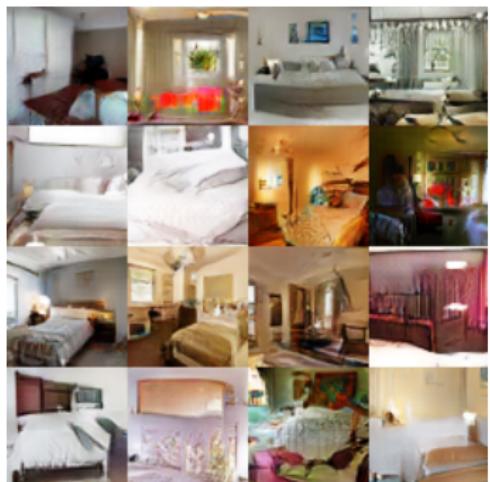
$\sim Q$

Task: training generative models

- Have: One collection of samples X from unknown distribution P .
- Goal: generate samples Q that look like P



LSUN bedroom samples P



Generated Q , MMD GAN

Using MMD to train a GAN

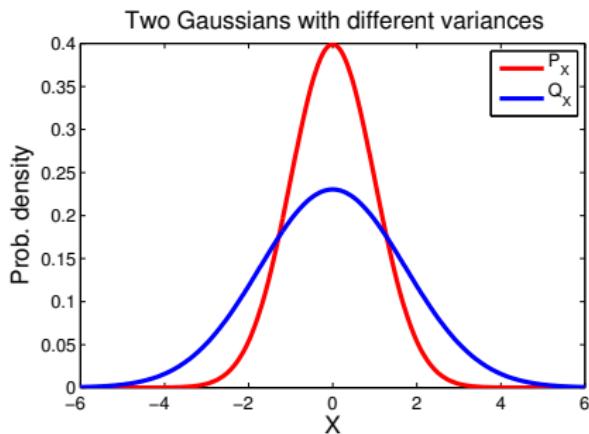
Outline

- Maximum Mean Discrepancy (MMD)...
 - ...as a difference in feature means
 - ...as an integral probability metric (not just a technicality!)
- Training GANs using MMD
 - Gradient penalty for training critic features
 - Comparison with WGAN-GP penalty
 - Evaluating GAN performance. Problems with Inception and FID.
 - Gradient bias for MMD GAN? For WGAN-GP? (not covered)

Maximum Mean Discrepancy

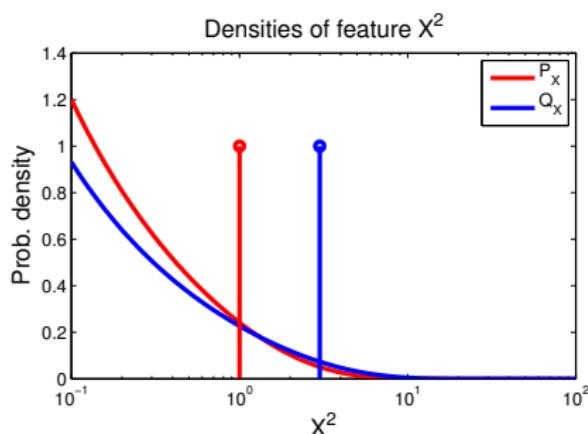
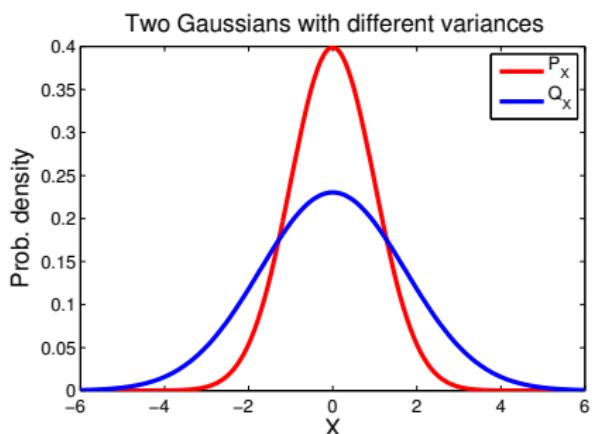
Feature mean difference

- Two Gaussians with same means, different variance
- Idea: look at difference in **means of features** of the RVs
- In Gaussian case: second order features of form $\varphi(x) = x^2$



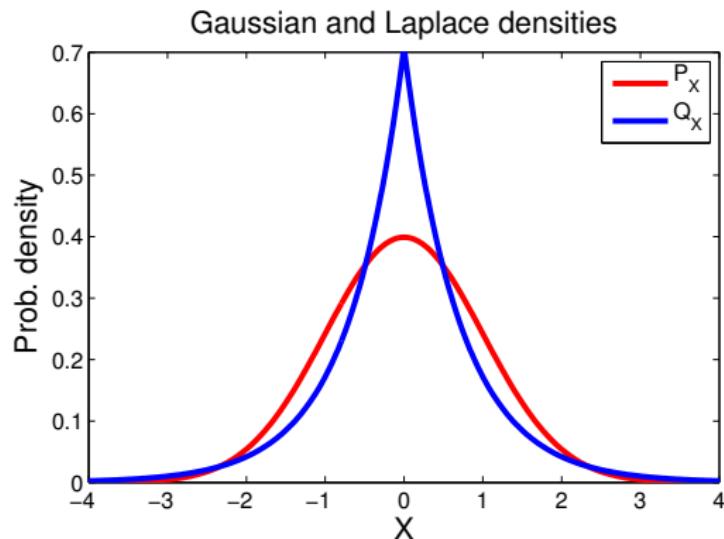
Feature mean difference

- Two Gaussians with same means, different variance
- Idea: look at difference in **means of features** of the RVs
- In Gaussian case: second order features of form $\varphi(x) = x^2$



Feature mean difference

- Gaussian and Laplace distributions
- Same mean and same variance
- Difference in means using **higher order features**...RKHS



Infinitely many features using kernels

Kernels: dot products
of features

Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_i(x) \dots] \in \ell_2$$

For positive definite k ,

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features
 $\varphi(x)$, dot product in
closed form!

Infinitely many features using kernels

Kernels: dot products
of features

Exponentiated quadratic kernel

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

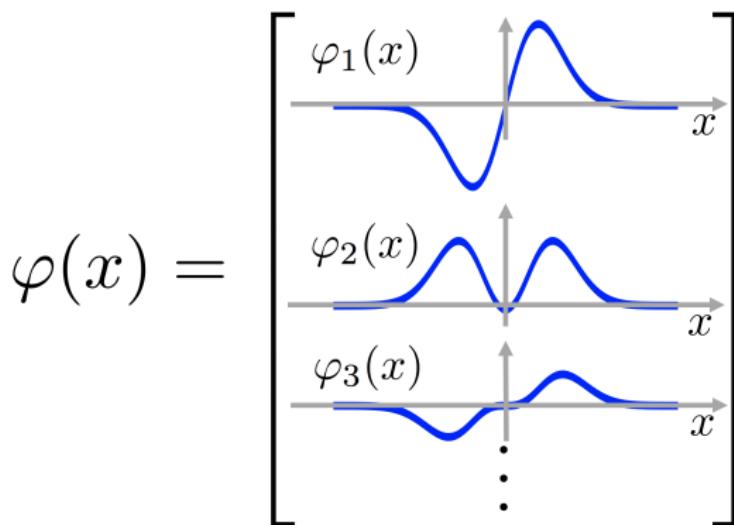
Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_i(x) \dots] \in \ell_2$$

For positive definite k ,

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features
 $\varphi(x)$, dot product in
closed form!



Infinitely many features of distributions

Given P a Borel **probability measure** on \mathcal{X} , define feature map of probability P ,

$$\mu_P = [\dots \mathbf{E}_P [\varphi_i(X)] \dots]$$

For positive definite $k(x, x')$,

$$\langle \mu_P, \mu_Q \rangle_{\mathcal{F}} = \mathbf{E}_{P, Q} k(\textcolor{teal}{x}, \textcolor{red}{y})$$

for $x \sim P$ and $y \sim Q$.

Fine print: feature map $\varphi(x)$ must be Bochner integrable for all probability measures considered.
Always true if kernel bounded.

Infinitely many features of distributions

Given P a Borel **probability measure** on \mathcal{X} , define feature map of probability P ,

$$\mu_P = [\dots \mathbf{E}_P [\varphi_i(X)] \dots]$$

For positive definite $k(x, x')$,

$$\langle \mu_P, \mu_Q \rangle_{\mathcal{F}} = \mathbf{E}_{P, Q} k(\textcolor{blue}{x}, \textcolor{red}{y})$$

for $x \sim P$ and $y \sim Q$.

Fine print: feature map $\varphi(x)$ must be Bochner integrable for all probability measures considered.
Always true if kernel bounded.

The maximum mean discrepancy

The maximum mean discrepancy is the distance between feature means:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P,Q} k(X, Y)}_{(b)} \end{aligned}$$

The maximum mean discrepancy

The maximum mean discrepancy is the distance between feature means:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P, Q} k(X, Y)}_{(b)} \end{aligned}$$

The maximum mean discrepancy

The maximum mean discrepancy is the distance between feature means:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P,Q} k(X, Y)}_{(b)} \end{aligned}$$

(a)= within distrib. similarity, (b)= cross-distrib. similarity.

Illustration of MMD

- Dogs ($= P$) and fish ($= Q$) example revisited
- Each entry is one of $k(\text{dog}_i, \text{dog}_j)$, $k(\text{dog}_i, \text{fish}_j)$, or $k(\text{fish}_i, \text{fish}_j)$

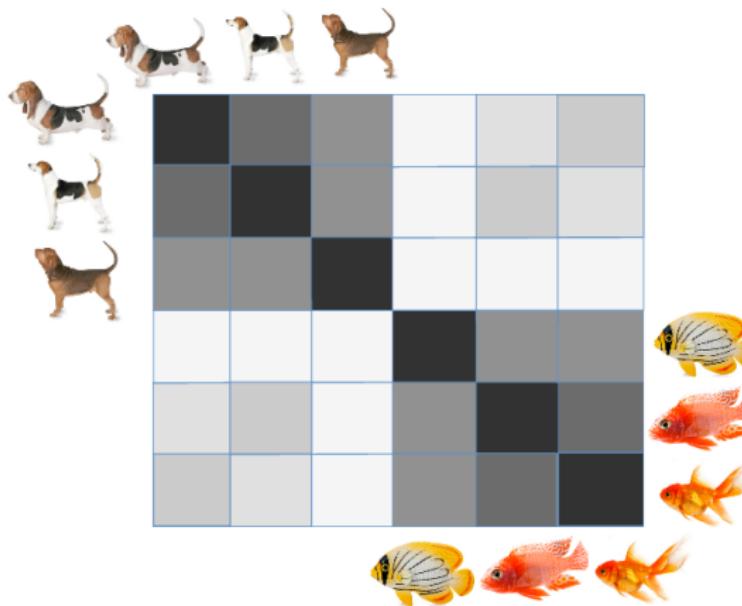
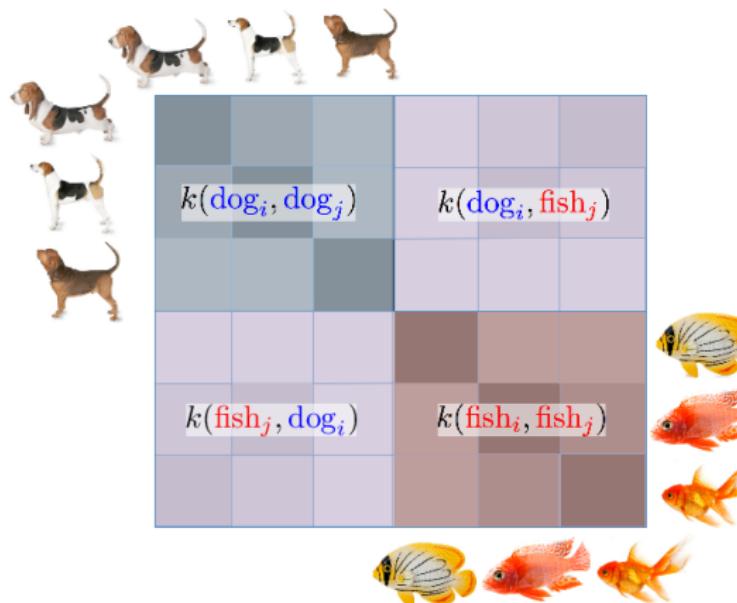


Illustration of MMD

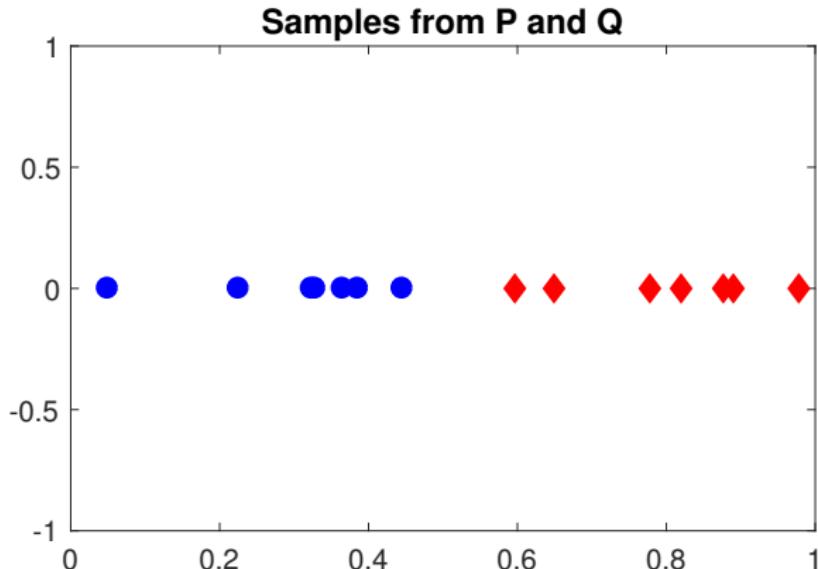
The maximum mean discrepancy:

$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{dog}_i, \text{dog}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{fish}_i, \text{fish}_j) - \frac{2}{n^2} \sum_{i,j} k(\text{dog}_i, \text{fish}_j)$$



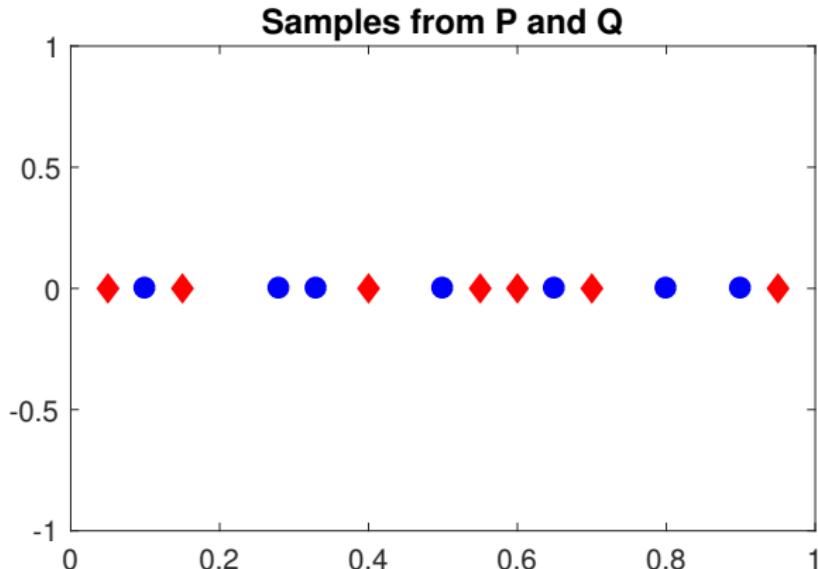
MMD as an integral probability metric

Are P and Q different?



MMD as an integral probability metric

Are P and Q different?

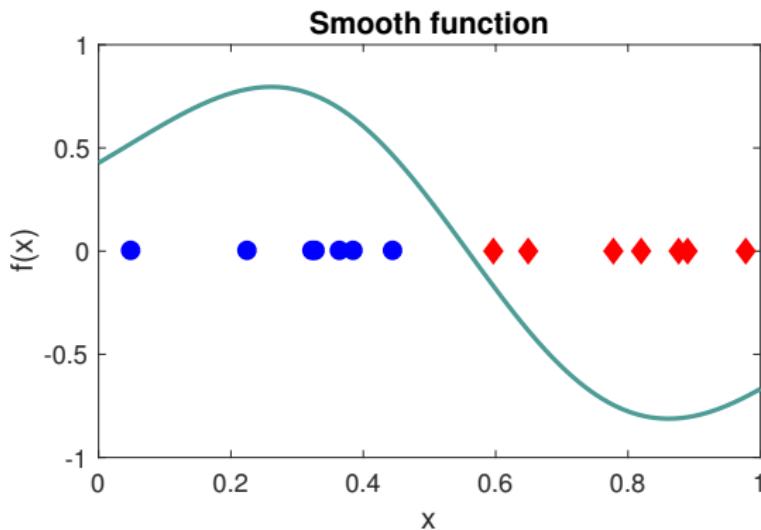


MMD as an integral probability metric

Integral probability metric:

Find a "well behaved function" $f(x)$ to maximize

$$\mathbf{E}_P f(\mathcal{X}) - \mathbf{E}_Q f(\mathcal{Y})$$

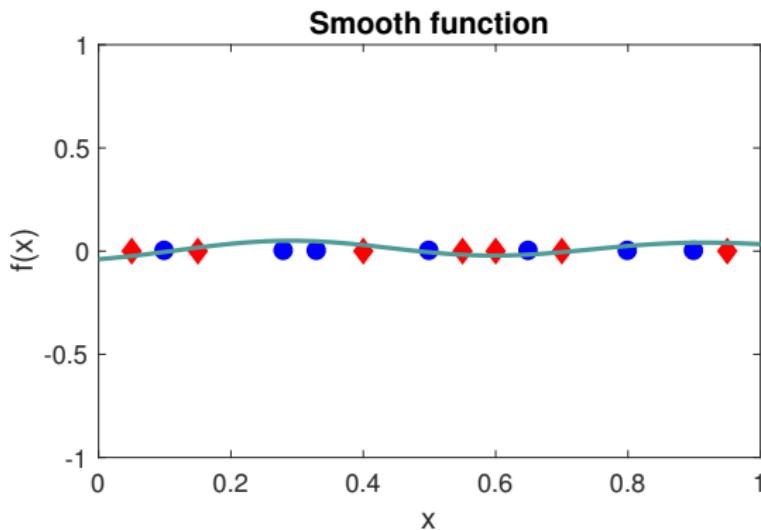


MMD as an integral probability metric

Integral probability metric:

Find a "well behaved function" $f(x)$ to maximize

$$\mathbf{E}_P f(\mathcal{X}) - \mathbf{E}_Q f(\mathcal{Y})$$



MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q; F) := \sup_{\|f\| \leq 1} [\mathbf{E}_{Pf}(X) - \mathbf{E}_{Qf}(Y)]$$

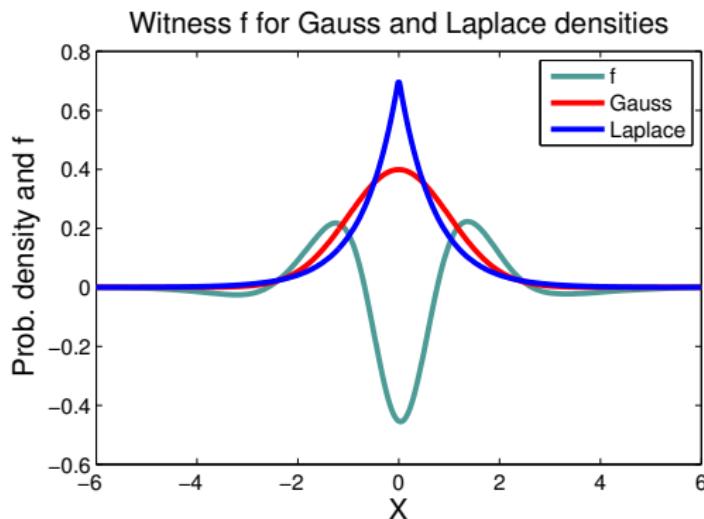
$(F = \text{unit ball in RKHS } \mathcal{F})$

MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q; F) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(F = \text{unit ball in RKHS } \mathcal{F})$



MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

Functions are linear combinations of features:

$$f(x) = \langle f, \varphi(x) \rangle_{\mathcal{F}} = \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \end{bmatrix}^T \begin{bmatrix} \varphi_1(x) \\ \varphi_2(x) \\ \varphi_3(x) \\ \vdots \end{bmatrix}$$

MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

Expectations of functions are linear combinations
of expected features

$$\mathbf{E}_P(f(X)) = \langle f, \mathbf{E}_P \varphi(X) \rangle_{\mathcal{F}} = \langle f, \mu_P \rangle_{\mathcal{F}}$$

(always true if kernel is bounded)

MMD as an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

For characteristic RKHS \mathcal{F} , $MMD(P, Q; \mathcal{F}) = 0$ iff $P = Q$

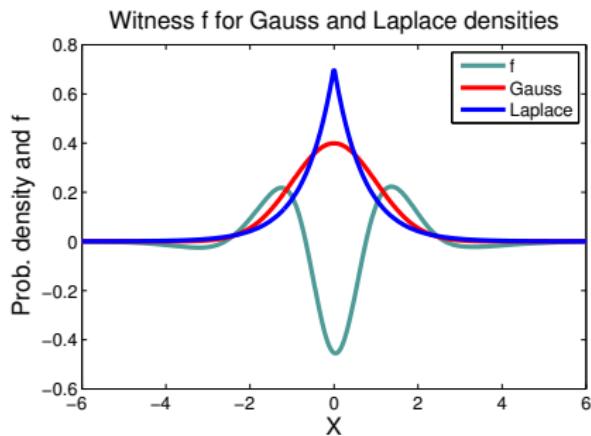
Other choices for witness function class:

- Bounded continuous [Dudley, 2002]
- Bounded variation 1 (Kolmogorov metric) [Müller, 1997]
- Bounded Lipschitz (Wasserstein distances) [Dudley, 2002]

Integral prob. metric vs feature difference

The MMD:

$$\begin{aligned} MMD(P, Q; F) \\ = \sup_{f \in F} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)] \end{aligned}$$



Integral prob. metric vs feature difference

The MMD:

use

$$\begin{aligned} MMD(\textcolor{blue}{P}, \textcolor{red}{Q}; \textcolor{teal}{F}) &= \sup_{f \in F} [\mathbf{E}_{\textcolor{blue}{P}} f(\textcolor{blue}{X}) - \mathbf{E}_{\textcolor{red}{Q}} f(\textcolor{red}{Y})] \\ &= \sup_{f \in F} \langle f, \boldsymbol{\mu}_P - \boldsymbol{\mu}_Q \rangle_{\mathcal{F}} \end{aligned}$$

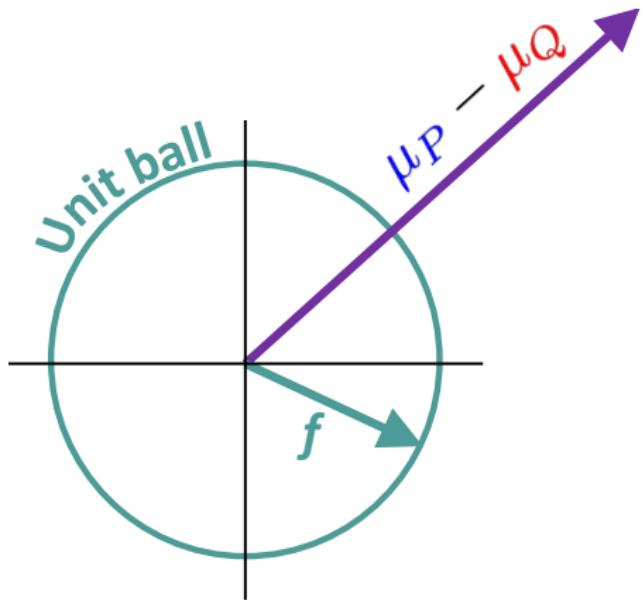
Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; F)$$

$$= \sup_{f \in F} [\mathbf{E}_P f(\mathbf{X}) - \mathbf{E}_Q f(\mathbf{Y})]$$

$$= \sup_{f \in F} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



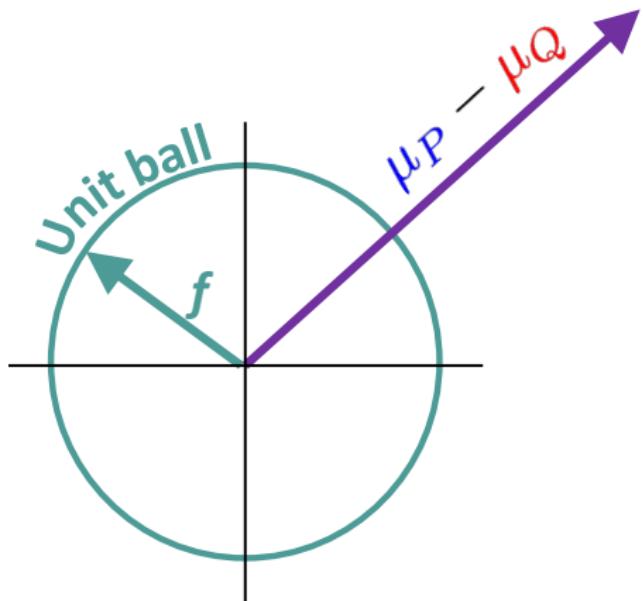
Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; F)$$

$$= \sup_{f \in F} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$$= \sup_{f \in F} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



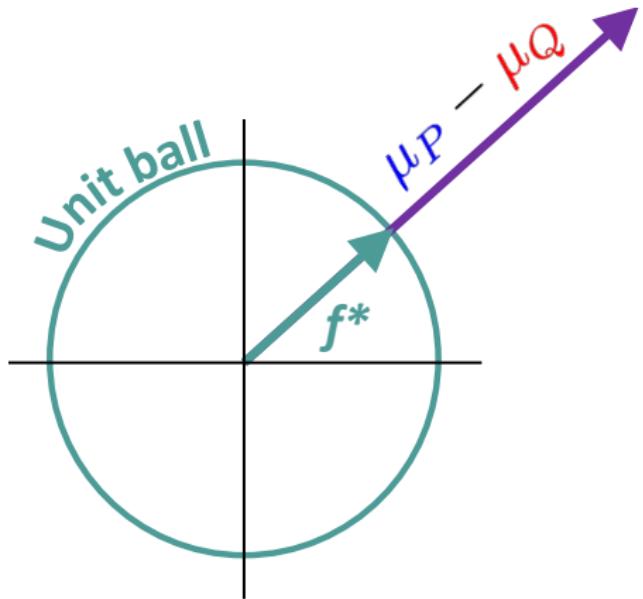
Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; F)$$

$$= \sup_{f \in F} [\mathbf{E}_P f(\mathbf{X}) - \mathbf{E}_Q f(\mathbf{Y})]$$

$$= \sup_{f \in F} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



$$f^* = \frac{\mu_P - \mu_Q}{\|\mu_P - \mu_Q\|}$$

Integral prob. metric vs feature difference

The MMD:

$$\begin{aligned}MMD(P, Q; F) &= \sup_{f \in F} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)] \\&= \sup_{f \in F} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}} \\&= \|\mu_P - \mu_Q\|\end{aligned}$$

Function view and feature view equivalent

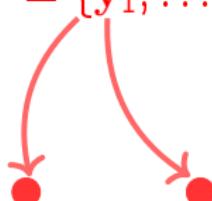
Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)

Observe $X = \{x_1, \dots, x_n\} \sim P$

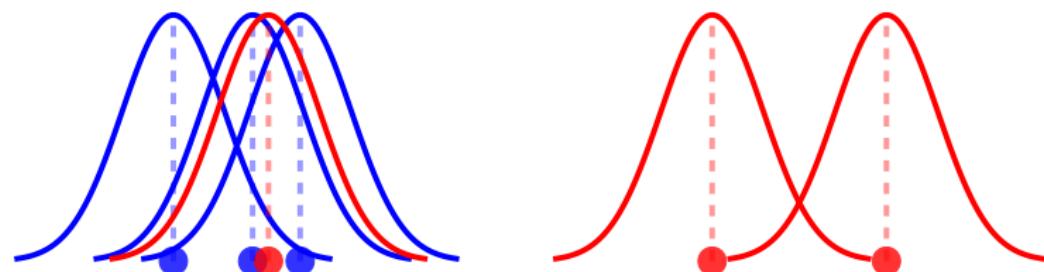


Observe $Y = \{y_1, \dots, y_n\} \sim Q$



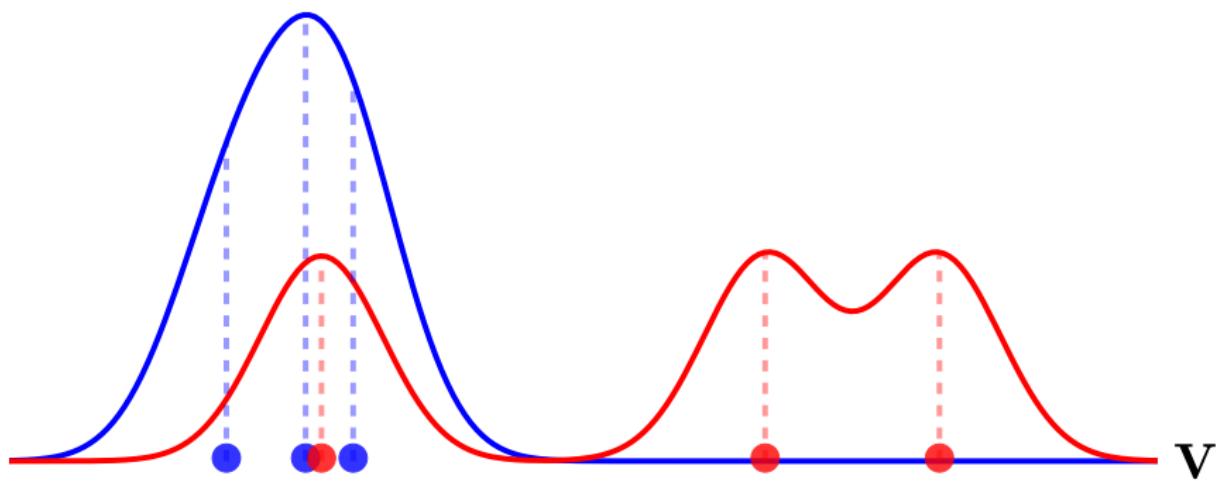
Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



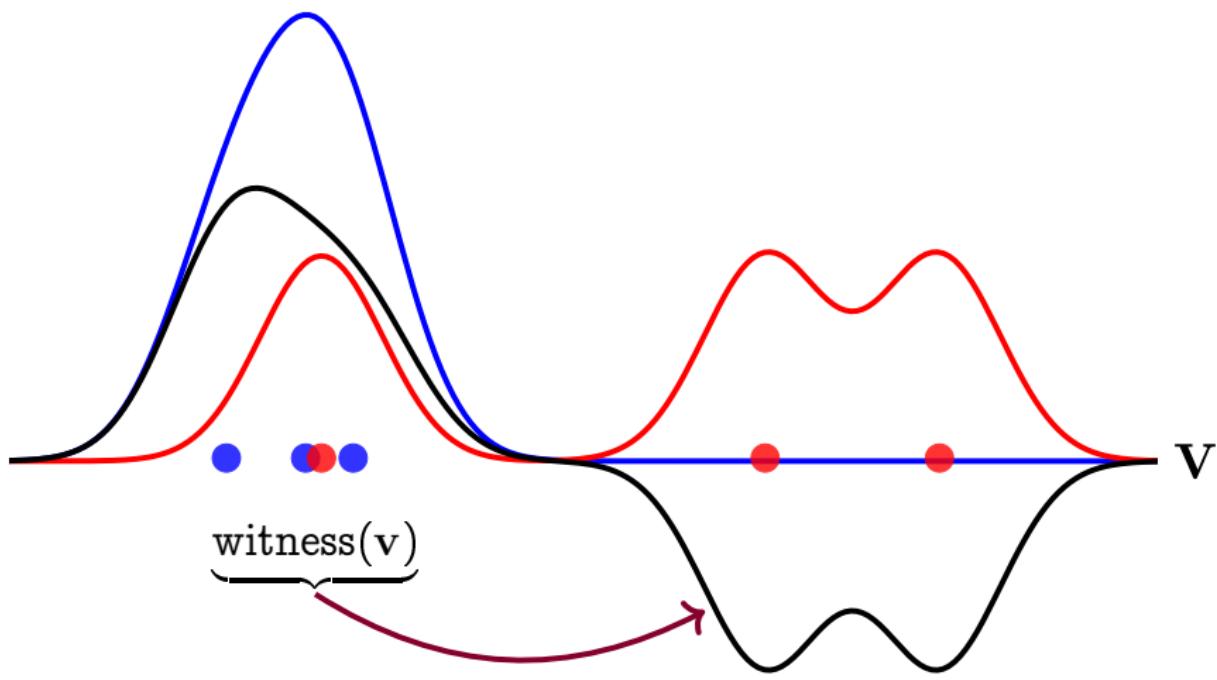
Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



Derivation of empirical witness function

Recall the witness function expression

$$f^* \propto \mu_P - \mu_Q$$

Derivation of empirical witness function

Recall the **witness function** expression

$$f^* \propto \mu_P - \mu_Q$$

The empirical feature mean for P

$$\hat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

Derivation of empirical witness function

Recall the **witness function** expression

$$\textcolor{teal}{f}^* \propto \mu_P - \mu_Q$$

The empirical feature mean for P

$$\widehat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

The empirical witness function at v

$$\textcolor{teal}{f}^*(v) = \langle \textcolor{teal}{f}^*, \varphi(v) \rangle_{\mathcal{F}}$$

Derivation of empirical witness function

Recall the **witness function** expression

$$f^* \propto \mu_P - \mu_Q$$

The empirical feature mean for P

$$\hat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

The empirical witness function at v

$$\begin{aligned} f^*(v) &= \langle f^*, \varphi(v) \rangle_{\mathcal{F}} \\ &\propto \langle \hat{\mu}_P - \hat{\mu}_Q, \varphi(v) \rangle_{\mathcal{F}} \end{aligned}$$

Derivation of empirical witness function

Recall the witness function expression

$$f^* \propto \mu_P - \mu_Q$$

The empirical feature mean for P

$$\hat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

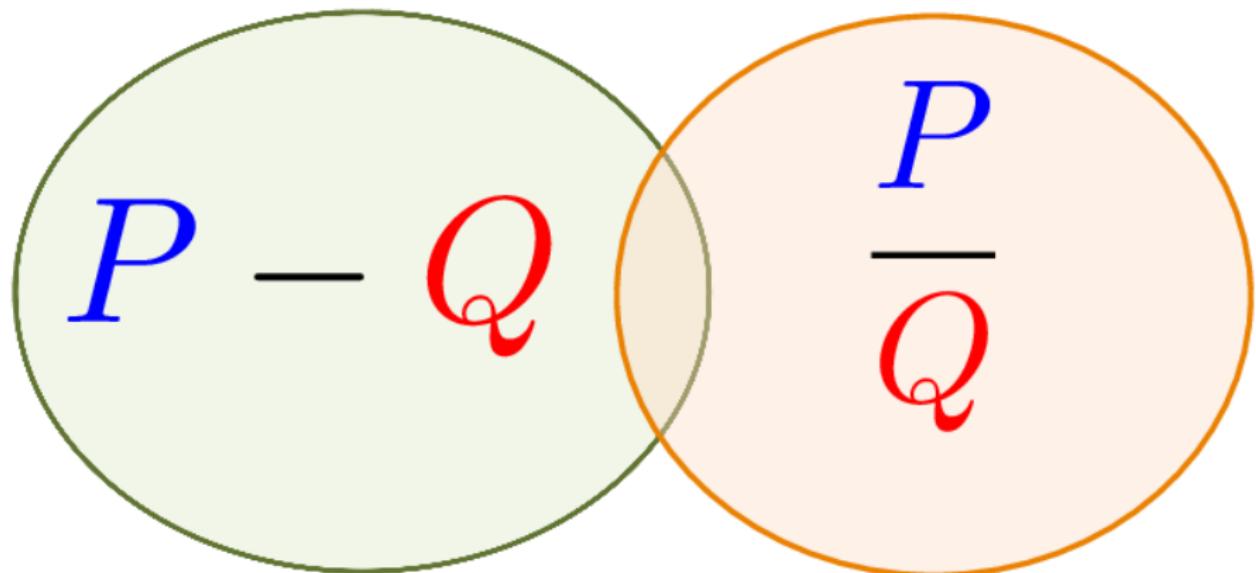
The empirical witness function at v

$$\begin{aligned} f^*(v) &= \langle f^*, \varphi(v) \rangle_{\mathcal{F}} \\ &\propto \langle \hat{\mu}_P - \hat{\mu}_Q, \varphi(v) \rangle_{\mathcal{F}} \\ &= \frac{1}{n} \sum_{i=1}^n k(\textcolor{blue}{x}_i, v) - \frac{1}{n} \sum_{i=1}^n k(\textcolor{red}{y}_i, v) \end{aligned}$$

Don't need explicit feature coefficients $f^* := [f_1^* \ f_2^* \ \dots]$

Interlude: divergence measures

Divergences



Divergences

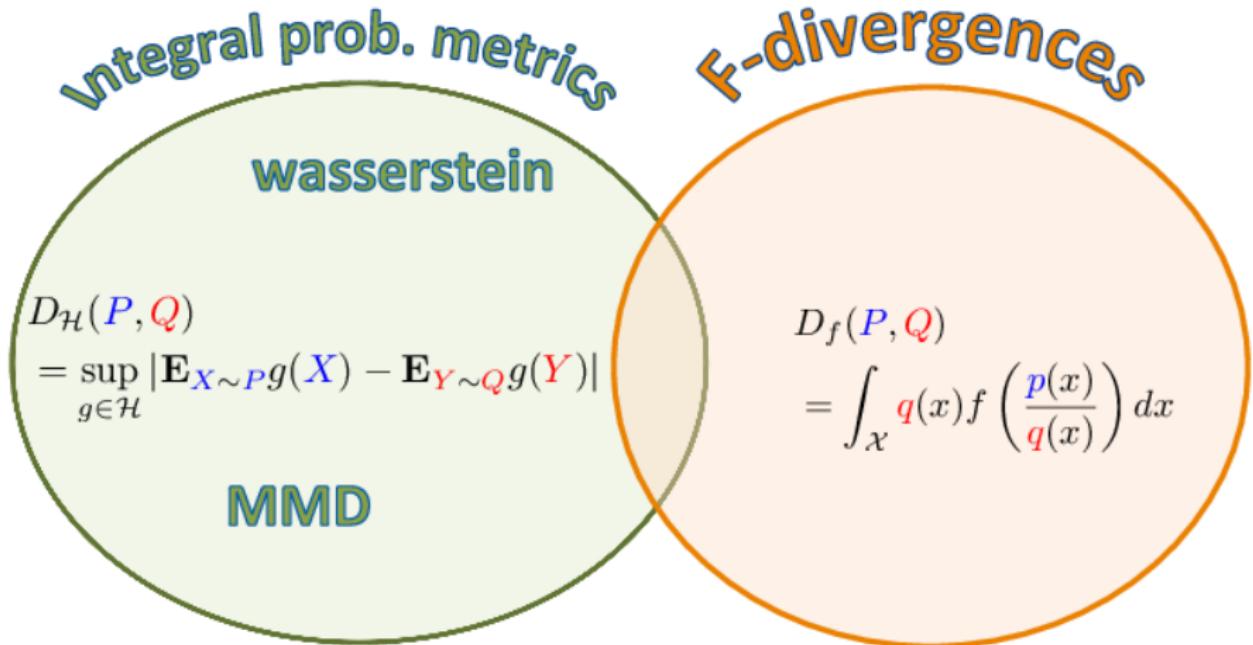
Integral prob. metrics

F-divergences

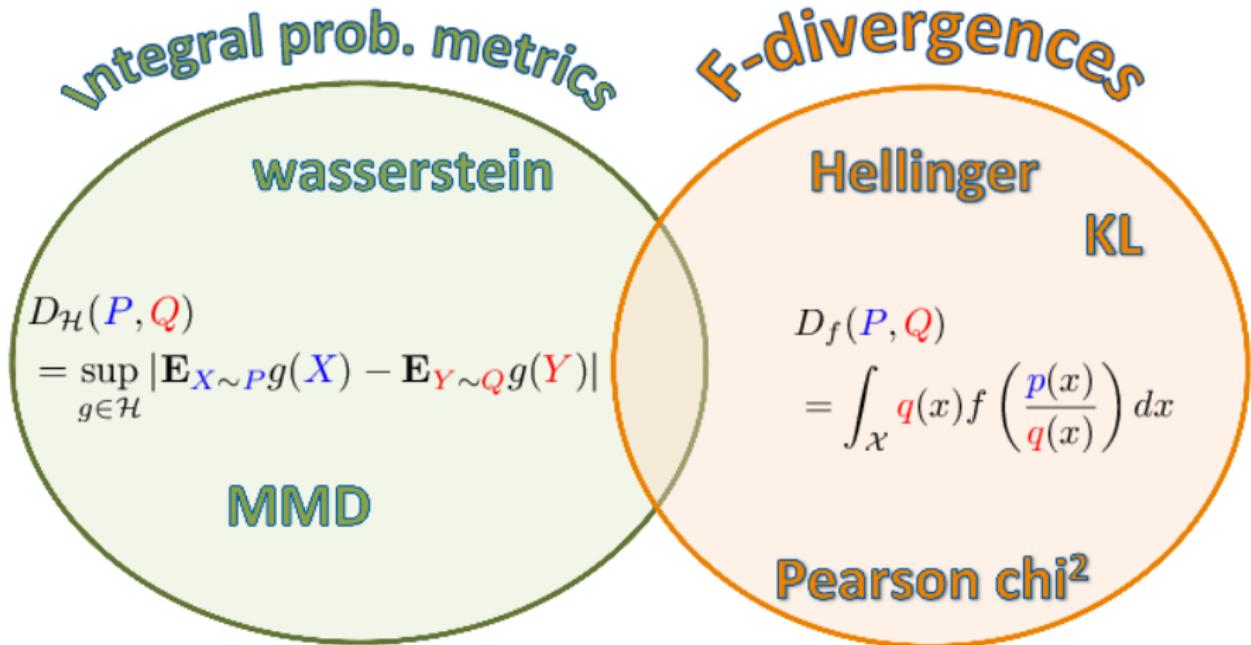
$$D_{\mathcal{H}}(\mathbf{P}, \mathbf{Q}) = \sup_{g \in \mathcal{H}} |\mathbf{E}_{X \sim \mathbf{P}} g(X) - \mathbf{E}_{Y \sim \mathbf{Q}} g(Y)|$$

$$D_f(\mathbf{P}, \mathbf{Q}) = \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

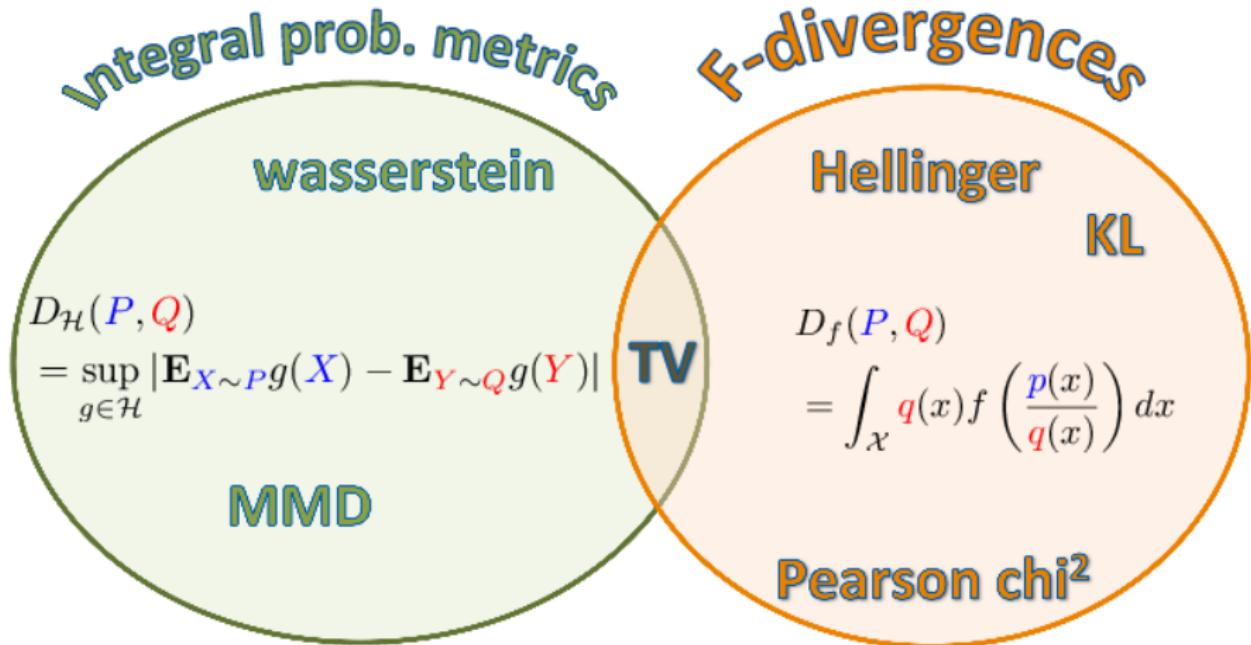
Divergences



Divergences



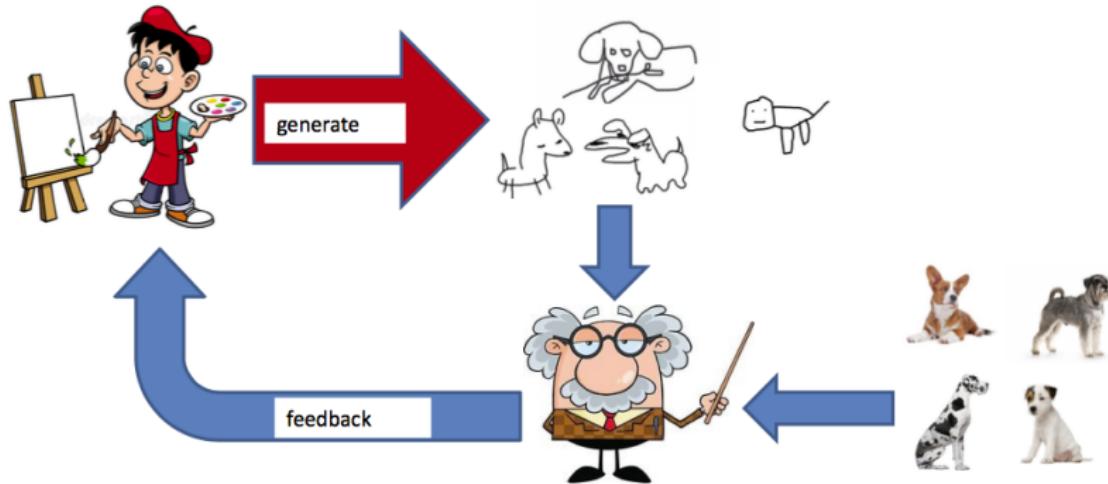
Divergences



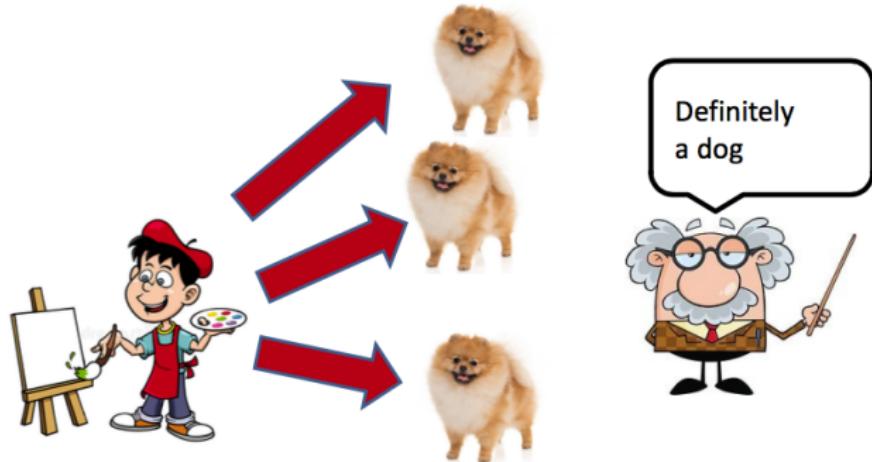
Sriperumbudur, Fukumizu, G, Schoelkopf, Lanckriet [EJS 2012]

Training GANs with MMD

Notation: Generative Adversarial Network (GAN)?



Why is classification not enough?



Classification **not** enough!
Need to compare **sets**

(otherwise student can just produce the **same dog** over and over)

MMD for GAN critic

Can you use MMD as a critic to train GANs?

From ICML 2015:

Generative Moment Matching Networks

Yujia Li¹

Kevin Swersky¹

Richard Zemel^{1,2}

YUJIALI@CS.TORONTO.EDU

KSWERSKY@CS.TORONTO.EDU

ZEMEL@CS.TORONTO.EDU

¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA

²Canadian Institute for Advanced Research, Toronto, ON, CANADA

From UAI 2015:

Training generative neural networks via Maximum Mean Discrepancy optimization

Gintare Karolina Dziugaite
University of Cambridge

Daniel M. Roy
University of Toronto

Zoubin Ghahramani
University of Cambridge

MMD for GAN critic

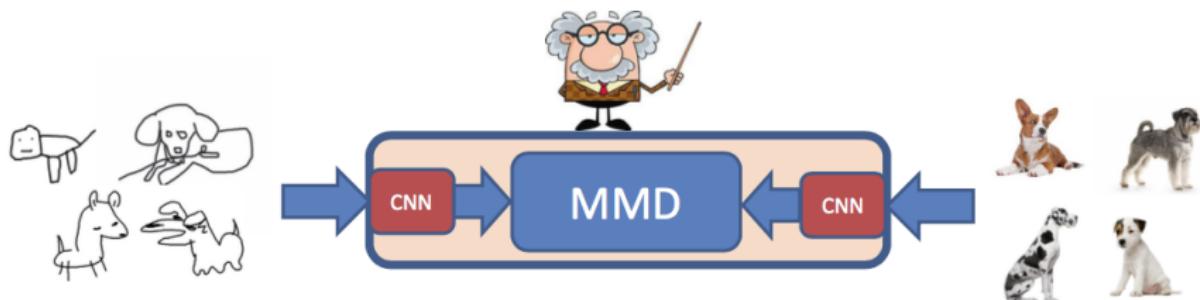
Can you use MMD as a critic to train GANs?



Need better image features.

How to improve the critic witness

- Add convolutional features!
- The **critic** (teacher) also needs to be trained.
- How to regularise?



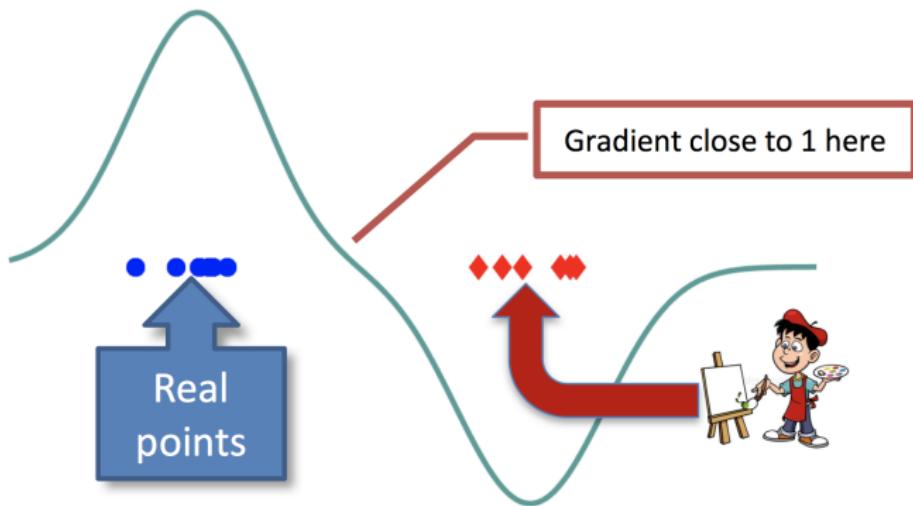
MMD GAN Li et al., [NIPS 2017]

Coulomb GAN Unterthiner et al., [ICLR 2018]

WGAN-GP

Wasserstein GAN Arjovsky et al. [ICML 2017].

WGAN-GP Gukrajani et al. [NIPS 2017]



WGAN-GP

Wasserstein GAN Arjovsky et al. [ICML 2017].

WGAN-GP Gukrajani et al. [NIPS 2017]



- Given a generator G_θ with parameters θ to be trained.
Samples $Y \sim G_\theta(Z)$ where $Z \sim R$



- Given critic features h_ψ with parameters ψ to be trained. f_ψ a linear function of h_ψ .

WGAN-GP

Wasserstein GAN Arjovsky et al. [ICML 2017].

WGAN-GP Gukrajani et al. [NIPS 2017]



- Given a generator G_θ with parameters θ to be trained.
Samples $\textcolor{red}{Y} \sim G_\theta(\textcolor{red}{Z})$ where $\textcolor{red}{Z} \sim \textcolor{red}{R}$



- Given critic features h_ψ with parameters ψ to be trained. f_ψ a linear function of h_ψ .

WGAN-GP gradient penalty:

$$\max_{\psi} \mathbf{E}_{X \sim \textcolor{blue}{P}} f_\psi(\textcolor{blue}{X}) - \mathbf{E}_{Z \sim \textcolor{red}{R}} f_\psi(G_\theta(\textcolor{red}{Z})) + \lambda \mathbf{E}_{\widetilde{X}} \left(\|\nabla_{\widetilde{X}} f_\theta(\widetilde{X})\| - 1 \right)^2$$

where

$$\widetilde{X} = \gamma \textcolor{blue}{x}_i + (1 - \gamma) G_\psi(\textcolor{red}{z}_j)$$

$$\gamma \sim \mathcal{U}([0, 1]) \quad x_i \in \{x_\ell\}_{\ell=1}^m \quad z_j \in \{z_\ell\}_{\ell=1}^n$$

The (W)MMD

Train MMD critic features with the witness function gradient penalty

Binkowski, Sutherland, Arbel, G. [ICLR 2018], Bellemare et al. [2017] for energy distance:

$$\max_{\psi} \text{MMD}^2(h_{\psi}(\mathbf{X}), h_{\psi}(G_{\theta}(\mathbf{Z}))) + \lambda \mathbf{E}_{\widetilde{\mathbf{X}}} \left(\|\nabla_{\widetilde{\mathbf{X}}} f_{\psi}(\widetilde{\mathbf{X}})\| - 1 \right)^2$$

where

$$f_{\psi}(\cdot) = \frac{1}{m} \sum_{i=1}^m k(h_{\psi}(\mathbf{x}_i), \cdot) - \frac{1}{n} \sum_{j=1}^n k(h_{\psi}(G_{\theta}(\mathbf{z}_j)), \cdot)$$


$$\widetilde{\mathbf{X}} = \gamma \mathbf{x}_i + (1 - \gamma) G_{\psi}(\mathbf{z}_j)$$

$$\gamma \sim \mathcal{U}([0, 1]) \quad \mathbf{x}_i \in \{\mathbf{x}_{\ell}\}_{\ell=1}^m \quad \mathbf{z}_j \in \{\mathbf{z}_{\ell}\}_{\ell=1}^n$$

Remark by Bottou et al. (2017): this modifies the function class. So critic is not an **33/53** in RKHS \mathcal{F} .

MMD for GAN critic: revisited

From ICLR 2018:

DEMYSTIFYING MMD GANs

Mikołaj Bińkowski*

Department of Mathematics

Imperial College London

mikbinkowski@gmail.com

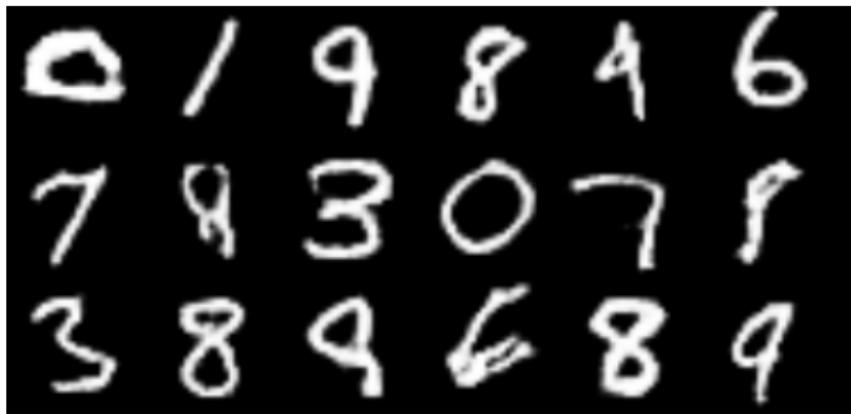
Dougal J. Sutherland,* Michael Arbel & Arthur Gretton

Gatsby Computational Neuroscience Unit

University College London

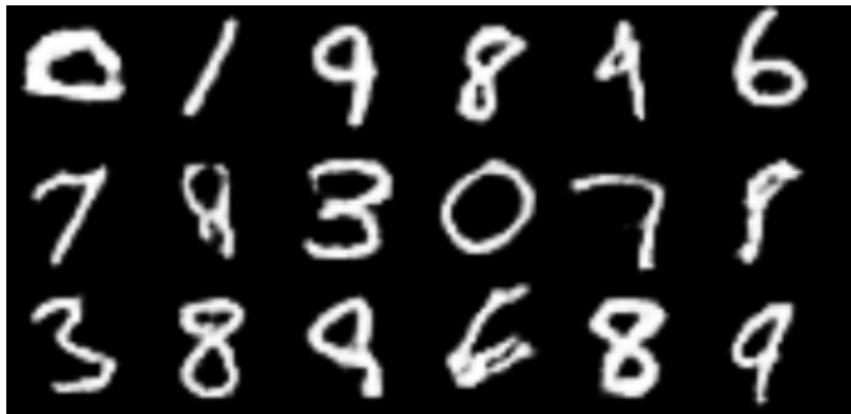
{dougal,michael.n.arbel,arthur.gretton}@gmail.com

MMD for GAN critic: revisited



Samples are better!

MMD for GAN critic: revisited



Samples are better!

Can we do better still?

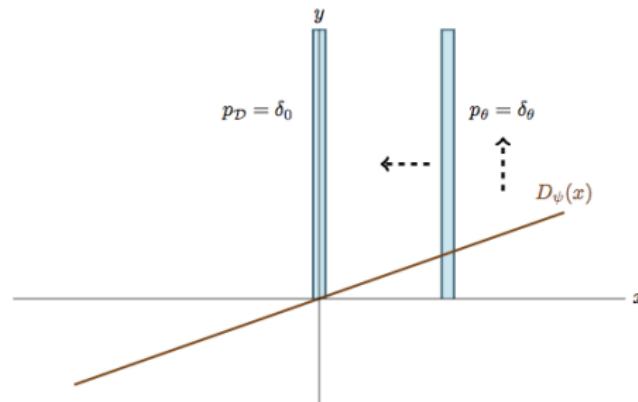
Convergence issues for WGAN-GP penalty

WGAN-GP style gradient penalty may not converge near solution

Nagarajan and Kolter [NIPS 2017], Mescheder et al. [ICML 2018], Balduzzi et al. [ICML 2018]

The Dirac-GAN

$$P = \delta_0 \quad Q = \delta_\theta \quad f_\psi(x) = \psi \cdot x$$



Convergence issues for WGAN-GP penalty

WGAN-GP style gradient penalty may not converge near solution

Nagarajan and Kolter [NIPS 2017], Mescheder et al. [ICML 2018], Balduzzi et al. [ICML 2018]

The Dirac-GAN

$$P = \delta_0 \quad Q = \delta_\theta \quad f_\psi(x) = \psi \cdot x$$

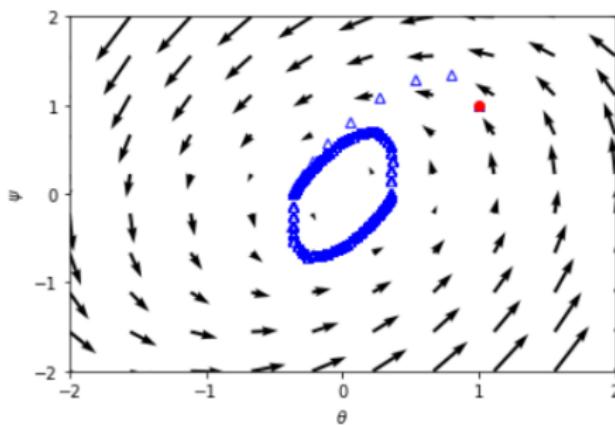


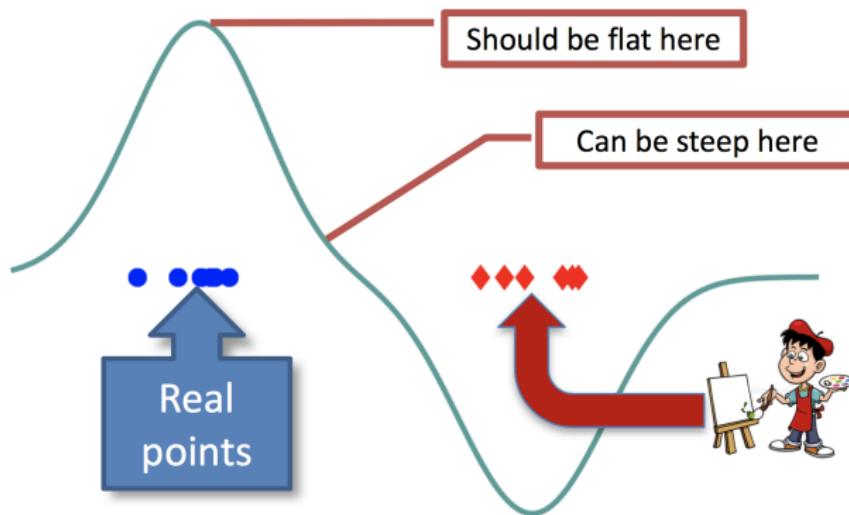
Figure from Mescheder et al. [ICML 2018]

A better gradient penalty

- New MMD GAN witness regulariser

Arbel, Sutherland, Binkowski, G. [arxiv, May 2018]

- Based on semi-supervised learning regulariser Bousquet et al. [NIPS 2004]
- Related to Sobolev GAN Mroueh et al. [ICLR 2018]



A better gradient penalty

■ New MMD GAN witness regulariser

Arbel, Sutherland, Binkowski, G. [arxiv, May 2018]

- Based on semi-supervised learning regulariser Bousquet et al. [NIPS 2004]
- Related to Sobolev GAN Mroueh et al. [ICLR 2018]

Modified witness function:

$$\widetilde{MMD} := \sup_{\|\mathbf{f}\|_S \leq 1} [\mathbb{E}_{\mathbf{P}} \mathbf{f}(\mathbf{X}) - \mathbb{E}_{\mathbf{Q}} \mathbf{f}(\mathbf{Y})]$$

where

$$\|\mathbf{f}\|_S^2 = \|\mathbf{f}\|_{L_2(\mathbf{P})}^2 + \|\nabla \mathbf{f}\|_{L_2(\mathbf{P})}^2 + \lambda \|\mathbf{f}\|_k^2$$

The equation shows the squared norm of \mathbf{f} in a space S as the sum of its squared L_2 norm over \mathbf{P} , the squared norm of its gradient over \mathbf{P} , and a term involving a regularization parameter λ and the squared norm of \mathbf{f} in a Reproducing Kernel Hilbert Space (RKHS) k . Three orange arrows point upwards from boxes labeled "L₂ norm control", "Gradient control", and "RKHS smoothness" to the corresponding terms in the equation.

A better gradient penalty

■ New MMD GAN witness regulariser

Arbel, Sutherland, Binkowski, G. [arxiv, May 2018]

■ Based on semi-supervised learning regulariser Bousquet et al. [NIPS 2004]

■ Related to Sobolev GAN Mroueh et al. [ICLR 2018]

Modified witness function:

$$\widetilde{MMD} := \sup_{\|\mathbf{f}\|_S \leq 1} [\mathbb{E}_{\mathbf{P}} \mathbf{f}(\mathbf{X}) - \mathbb{E}_{\mathbf{Q}} \mathbf{f}(\mathbf{Y})]$$

where

$$\|\mathbf{f}\|_S^2 = \|\mathbf{f}\|_{L_2(\mathbf{P})}^2 + \|\nabla \mathbf{f}\|_{L_2(\mathbf{P})}^2 + \lambda \|\mathbf{f}\|_k^2$$

The equation shows the squared norm of \mathbf{f} as the sum of three terms: the squared L_2 norm of \mathbf{f} over the domain \mathbf{P} , the squared gradient norm of \mathbf{f} over the same domain, and a regularization term involving the k -norm of \mathbf{f} . Below the equation, three orange arrows point upwards from three boxes to their respective terms: 'L₂ norm control' points to the first term, 'Gradient control' points to the second, and 'RKHS smoothness' points to the third.

Problem: not computationally feasible: $O(n^3)$ per iteration.

A better gradient penalty

■ New MMD GAN witness regulariser

Arbel, Sutherland, Binkowski, G. [arxiv, May 2018]

- Based on semi-supervised learning regulariser Bousquet et al. [NIPS 2004]
- Related to Sobolev GAN Mroueh et al. [ICLR 2018]

The scaled MMD:

$$SMMD = \sigma_{k, P, \lambda} MMD$$

where

$$\sigma_{k, P, \lambda} = \left(\lambda + \int k(x, x) dP(x) + \sum_{i=1}^d \int \partial_i \partial_{i+d} k(x, y) dP(x) \right)^{-1/2}$$

Replace expensive constraint with cheap upper bound:

$$\|\mathbf{f}\|_S^2 \leq \sigma_{k, P, \lambda}^{-1} \|\mathbf{f}\|_k^2$$

A better gradient penalty

■ New MMD GAN witness regulariser

Arbel, Sutherland, Binkowski, G. [arxiv, May 2018]

- Based on semi-supervised learning regulariser Bousquet et al. [NIPS 2004]
- Related to Sobolev GAN Mroueh et al. [ICLR 2018]

The scaled MMD:

$$SMMD = \sigma_{k, P, \lambda} MMD$$

where

$$\sigma_{k, P, \lambda} = \left(\lambda + \int k(x, x) dP(x) + \sum_{i=1}^d \int \partial_i \partial_{i+d} k(x, y) dP(x) \right)^{-1/2}$$

Replace expensive constraint with cheap upper bound:

$$\|\mathbf{f}\|_S^2 \leq \sigma_{k, P, \lambda}^{-1} \|\mathbf{f}\|_k^2$$

Idea: rather than regularise the critic or witness function, regularise features directly

Evaluation and experiments

Evaluation of GANs

The inception score? Salimans et al. [NIPS 2016]

Based on the classification output $p(y|x)$ of the inception model Szegedy et al. [ICLR 2014],

$$E_x \exp KL(P(y|x) \| P(y)).$$

High when:

- predictive label distribution $P(y|x)$ has low entropy (good quality images)
- label entropy $P(y)$ is high (good variety).

Evaluation of GANs

The inception score? Salimans et al. [NIPS 2016]

Based on the classification output $p(y|x)$ of the inception model Szegedy et al. [ICLR 2014],

$$E_x \exp KL(P(y|x) \| P(y)).$$

High when:

- predictive label distribution $P(y|x)$ has low entropy (good quality images)
- label entropy $P(y)$ is high (good variety).

Problem: relies on a trained classifier! Can't be used on new categories (celeb, bedroom...)

Evaluation of GANs

The Frechet inception distance? Heusel et al. [NIPS 2017]

Fits Gaussians to features in the inception architecture (pool3 layer):

$$FID(\mathcal{P}, \mathcal{Q}) = \|\mu_{\mathcal{P}} - \mu_{\mathcal{Q}}\|^2 + \text{tr}(\Sigma_{\mathcal{P}}) + \text{tr}(\Sigma_{\mathcal{Q}}) - 2\text{tr}\left((\Sigma_{\mathcal{P}}\Sigma_{\mathcal{Q}})^{\frac{1}{2}}\right)$$

where $\mu_{\mathcal{P}}$ and $\Sigma_{\mathcal{P}}$ are the feature mean and covariance of \mathcal{P}

Evaluation of GANs

The Frechet inception distance? Heusel et al. [NIPS 2017]

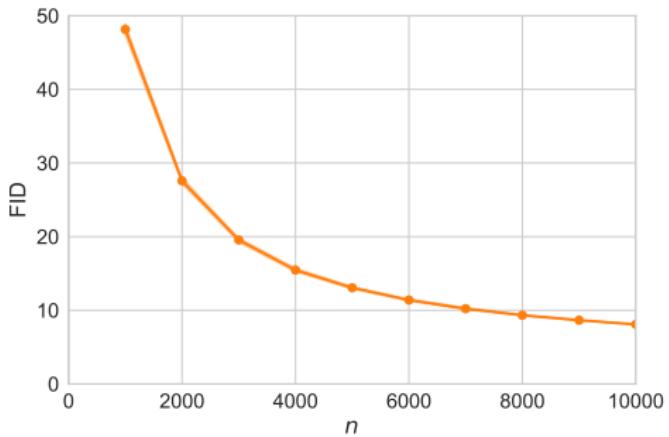
Fits Gaussians to features in the inception architecture (pool3 layer):

$$FID(P, Q) = \|\mu_P - \mu_Q\|^2 + \text{tr}(\Sigma_P) + \text{tr}(\Sigma_Q) - 2\text{tr}\left((\Sigma_P \Sigma_Q)^{\frac{1}{2}}\right)$$

where μ_P and Σ_P are the feature mean and covariance of P

Problem: bias. For finite samples can consistently give incorrect answer.

- Bias demo,
CIFAR-10 train vs
test



Evaluation of GANs

The FID can give the wrong answer in theory.

Assume m samples from P and $n \rightarrow \infty$ samples from Q .

Given two alternatives:

$$P_1 \sim \mathcal{N}(0, (1 - m^{-1})^2) \quad P_2 \sim \mathcal{N}(0, 1) \quad Q \sim \mathcal{N}(0, 1).$$

Clearly,

$$FID(P_1, Q) = \frac{1}{m^2} > FID(P_2, Q) = 0$$

Given m samples from P_1 and P_2 ,

$$FID(\widehat{P}_1, Q) < FID(\widehat{P}_2, Q).$$

Evaluation of GANs

The FID can give the wrong answer in theory.

Assume m samples from P and $n \rightarrow \infty$ samples from Q .

Given two alternatives:

$$P_1 \sim \mathcal{N}(0, (1 - m^{-1})^2) \quad P_2 \sim \mathcal{N}(0, 1) \quad Q \sim \mathcal{N}(0, 1).$$

Clearly,

$$FID(P_1, Q) = \frac{1}{m^2} > FID(P_2, Q) = 0$$

Given m samples from P_1 and P_2 ,

$$FID(\widehat{P}_1, Q) < FID(\widehat{P}_2, Q).$$

Evaluation of GANs

The FID can give the wrong answer in theory.

Assume m samples from P and $n \rightarrow \infty$ samples from Q .

Given two alternatives:

$$P_1 \sim \mathcal{N}(0, (1 - m^{-1})^2) \quad P_2 \sim \mathcal{N}(0, 1) \quad Q \sim \mathcal{N}(0, 1).$$

Clearly,

$$FID(P_1, Q) = \frac{1}{m^2} > FID(P_2, Q) = 0$$

Given m samples from P_1 and P_2 ,

$$FID(\widehat{P}_1, Q) < FID(\widehat{P}_2, Q).$$

Evaluation of GANs

The FID can give the wrong answer in theory.

Assume m samples from P and $n \rightarrow \infty$ samples from Q .

Given two alternatives:

$$P_1 \sim \mathcal{N}(0, (1 - m^{-1})^2) \quad P_2 \sim \mathcal{N}(0, 1) \quad Q \sim \mathcal{N}(0, 1).$$

Clearly,

$$FID(P_1, Q) = \frac{1}{m^2} > FID(P_2, Q) = 0$$

Given m samples from P_1 and P_2 ,

$$FID(\widehat{P_1}, Q) < FID(\widehat{P_2}, Q).$$

Evaluation of GANs

The FID can give the wrong answer in practice.

Let $d = 2048$, and define

$$P_1 = \text{relu}(\mathcal{N}(0, I_d)) \quad P_2 = \text{relu}(\mathcal{N}(1, .8\Sigma + .2I_d)) \quad Q = \text{relu}(\mathcal{N}(1, I_d))$$

where $\Sigma = \frac{4}{d}CC^T$, with C a $d \times d$ matrix with iid standard normal entries.

For a random draw of C :

$$FID(P_1, Q) \approx 1123.0 > 1114.8 \approx FID(P_2, Q)$$

With $m = 50\,000$ samples,

$$FID(\widehat{P}_1, Q) \approx 1133.7 < 1136.2 \approx FID(\widehat{P}_2, Q)$$

At $m = 100\,000$ samples, the ordering of the estimates is correct.

This behavior is similar for other random draws of C .

Evaluation of GANs

The FID can give the wrong answer in practice.

Let $d = 2048$, and define

$$\textcolor{blue}{P}_1 = \text{relu}(\mathcal{N}(\mathbf{0}, I_d)) \quad \textcolor{blue}{P}_2 = \text{relu}(\mathcal{N}(\mathbf{1}, .8\Sigma + .2I_d)) \quad \textcolor{red}{Q} = \text{relu}(\mathcal{N}(\mathbf{1}, I_d))$$

where $\Sigma = \frac{4}{d}CC^T$, with C a $d \times d$ matrix with iid standard normal entries.

For a random draw of C :

$$FID(\textcolor{blue}{P}_1, \textcolor{red}{Q}) \approx 1123.0 > 1114.8 \approx FID(\textcolor{blue}{P}_2, \textcolor{red}{Q})$$

With $m = 50\,000$ samples,

$$FID(\widehat{\textcolor{blue}{P}}_1, \textcolor{red}{Q}) \approx 1133.7 < 1136.2 \approx FID(\widehat{\textcolor{blue}{P}}_2, \textcolor{red}{Q})$$

At $m = 100\,000$ samples, the ordering of the estimates is correct.

This behavior is similar for other random draws of C .

Evaluation of GANs

The FID can give the wrong answer in practice.

Let $d = 2048$, and define

$$\textcolor{blue}{P}_1 = \text{relu}(\mathcal{N}(\mathbf{0}, I_d)) \quad \textcolor{blue}{P}_2 = \text{relu}(\mathcal{N}(\mathbf{1}, .8\Sigma + .2I_d)) \quad \textcolor{red}{Q} = \text{relu}(\mathcal{N}(\mathbf{1}, I_d))$$

where $\Sigma = \frac{4}{d}CC^T$, with C a $d \times d$ matrix with iid standard normal entries.

For a random draw of C :

$$FID(\textcolor{blue}{P}_1, \textcolor{red}{Q}) \approx 1123.0 > 1114.8 \approx FID(\textcolor{blue}{P}_2, \textcolor{red}{Q})$$

With $m = 50\,000$ samples,

$$FID(\widehat{\textcolor{blue}{P}_1}, \textcolor{red}{Q}) \approx 1133.7 < 1136.2 \approx FID(\widehat{\textcolor{blue}{P}_2}, \textcolor{red}{Q})$$

At $m = 100\,000$ samples, the ordering of the estimates is correct.

This behavior is similar for other random draws of C .

Evaluation of GANs

The FID can give the wrong answer in practice.

Let $d = 2048$, and define

$$\textcolor{blue}{P}_1 = \text{relu}(\mathcal{N}(\mathbf{0}, I_d)) \quad \textcolor{blue}{P}_2 = \text{relu}(\mathcal{N}(\mathbf{1}, .8\Sigma + .2I_d)) \quad \textcolor{red}{Q} = \text{relu}(\mathcal{N}(\mathbf{1}, I_d))$$

where $\Sigma = \frac{4}{d}CC^T$, with C a $d \times d$ matrix with iid standard normal entries.

For a random draw of C :

$$FID(\textcolor{blue}{P}_1, \textcolor{red}{Q}) \approx 1123.0 > 1114.8 \approx FID(\textcolor{blue}{P}_2, \textcolor{red}{Q})$$

With $m = 50\,000$ samples,

$$FID(\widehat{\textcolor{blue}{P}_1}, \textcolor{red}{Q}) \approx 1133.7 < 1136.2 \approx FID(\widehat{\textcolor{blue}{P}_2}, \textcolor{red}{Q})$$

At $m = 100\,000$ samples, the ordering of the estimates is correct.

This behavior is similar for other random draws of C .

The kernel inception distance (KID)

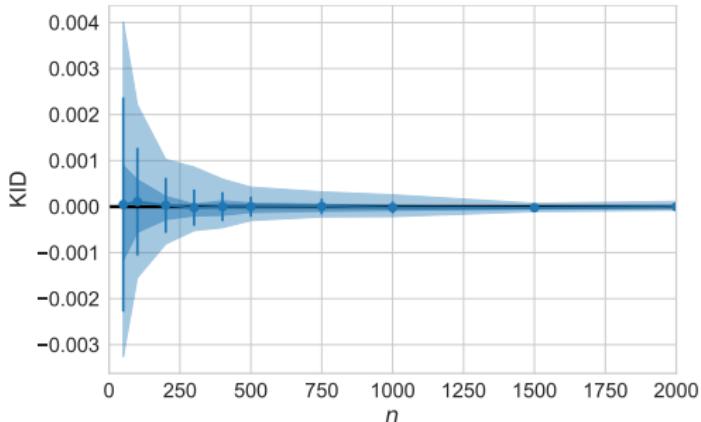
The Kernel inception distance Binkowski, Sutherland, Arbel, G. [ICLR 2018]

Measures similarity of the samples' representations in the inception architecture (pool3 layer)

MMD with kernel

$$k(x, y) = \left(\frac{1}{d} x^\top y + 1 \right)^3.$$

- Checks match for feature means, variances, skewness
- **Unbiased** : eg CIFAR-10 train/test



The kernel inception distance (KID)

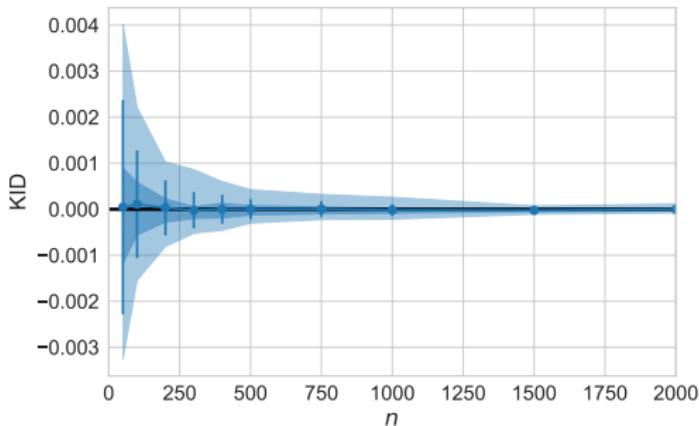
The Kernel inception distance Binkowski, Sutherland, Arbel, G. [ICLR 2018]

Measures similarity of the samples' representations in the inception architecture (pool3 layer)

MMD with kernel

$$k(x, y) = \left(\frac{1}{d} x^\top y + 1 \right)^3.$$

- Checks match for feature means, variances, skewness
- **Unbiased** : eg CIFAR-10 train/test



...“but isn't KID computationally costly?”

The kernel inception distance (KID)

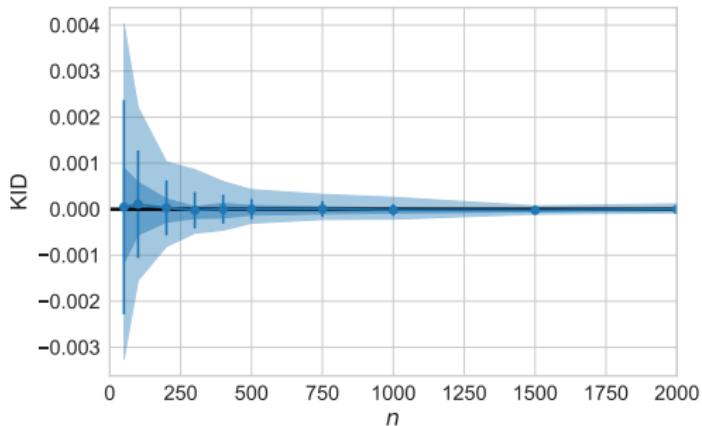
The Kernel inception distance Binkowski, Sutherland, Arbel, G. [ICLR 2018]

Measures similarity of the samples' representations in the inception architecture (pool3 layer)

MMD with kernel

$$k(x, y) = \left(\frac{1}{d} x^\top y + 1 \right)^3.$$

- Checks match for feature means, variances, skewness
- **Unbiased** : eg CIFAR-10 train/test



...“but isn't KID is computationally costly?”

“Block” KID implementation is cheaper than FID: see paper (or use our code)!

The kernel inception distance (KID)

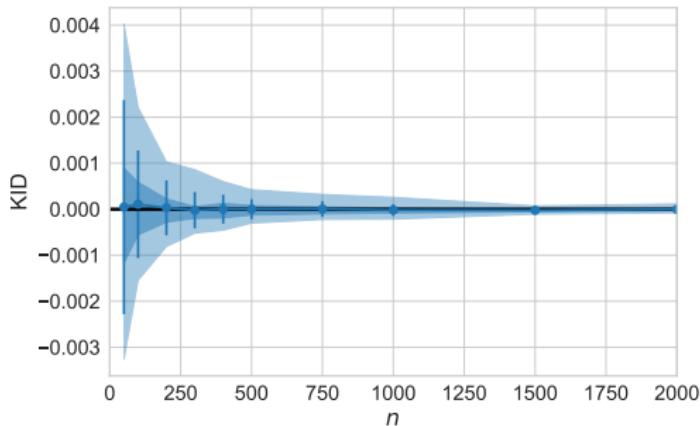
The Kernel inception distance Binkowski, Sutherland, Arbel, G. [ICLR 2018]

Measures similarity of the samples' representations in the inception architecture (pool3 layer)

MMD with kernel

$$k(x, y) = \left(\frac{1}{d} x^\top y + 1 \right)^3.$$

- Checks match for feature means, variances, skewness
- Unbiased : eg CIFAR-10 train/test



Also used for automatic learning rate adjustment: if $KID(\hat{P}_{t+1}, Q)$ not significantly better than $KID(\hat{P}_t, Q)$ then reduce learning rate.

[Bounliphone et al. ICLR 2016]

Benchmarks for comparison (all from ICLR 2018)

SPECTRAL NORMALIZATION FOR GENERATIVE ADVERSARIAL NETWORKS

Takeru Miyato¹, Toshiki Kataoka¹, Masanori Koyama², Yuichi Yoshida³

{miyato, kataoka}@preferred.jp

toyama.masanori@gmail.com

yoshi@li.ac.jp

¹Preferred Networks, Inc. ²Ritsumeikan University ³National Institute of Informatics

We combine with scaled MMD

DEMYSTIFYING MMD GANS

Mikolaj Binkowski*

Department of Mathematics

Imperial College London

mikbinkowski@gmail.com

Dougal J. Sutherland*, Michael Arbel & Arthur Gretton

Gatsby Computational Neuroscience Unit

University College London

dougal.sutherland, michael.n.arbel, arthur.gretton@gmail.com

Our ICLR
2018
paper

SOBOLEV GAN

Youssef Mroueh¹, Chun-Liang Li^{2,*}, Tom Sercombe^{1,*}, Anant Raj^{3,*} & Yu Cheng¹

† IBM Research AI

◦ Carnegie Mellon University

◊ Max Planck Institute for Intelligent Systems

* denotes Equal Contribution

{mroueh, chengyu}@us.ibm.com, chunliail@cs.cmu.edu,

tom.sercombe@ibm.com, anant.raj@tuebingen.mpg.de

BOUNDARY-SEEKING GENERATIVE ADVERSARIAL NETWORKS

R Devon Hjelm*

MILA, University of Montréal, IVADO

erroneous@gmail.com

Athul Paul Jacob*

MILA, MSR, University of Waterloo

apjacob@edu.uwaterloo.ca

Tong Che

MILA, University of Montréal

tong.che@umontreal.ca

Adam Trischler

MSR

adam.trischler@microsoft.com

Kyunghyun Cho

New York University

CIFAR Azrieli Global Scholar

kyunghyun.cho@nyu.edu

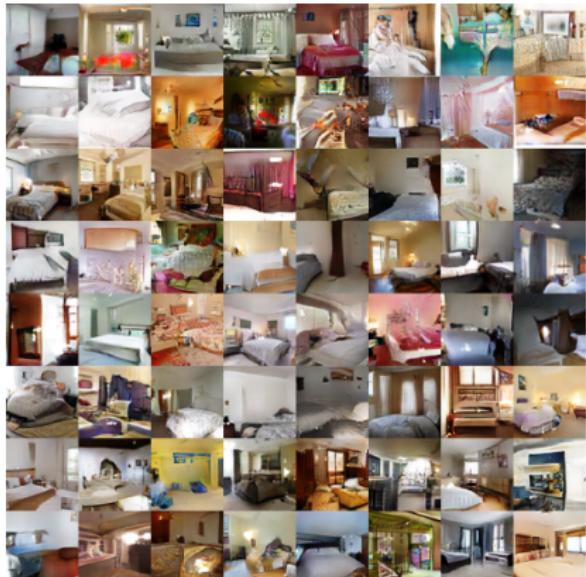
Yoshua Bengio

MILA, University of Montréal, CIFAR, IVADO

yoshua.bengio@umontreal.ca

Results: what does MMD buy you?

- Critic features from DCGAN: an f -filter critic has f , $2f$, $4f$ and $8f$ convolutional filters in layers 1-4. LSUN 64×64 .



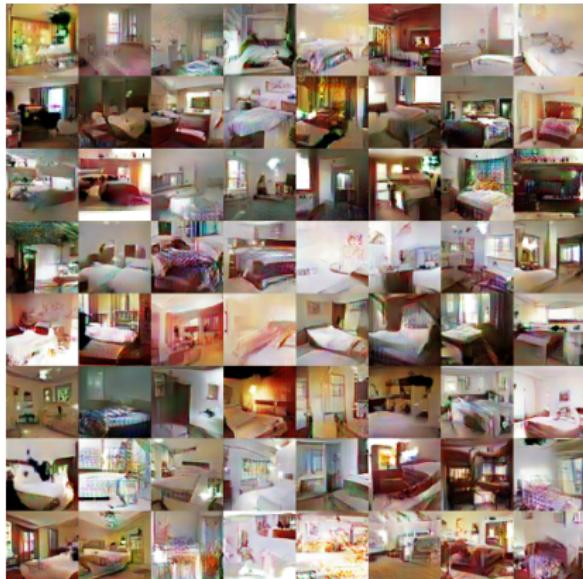
MMD GAN samples, $f = 64$,
FID=32, KID=3



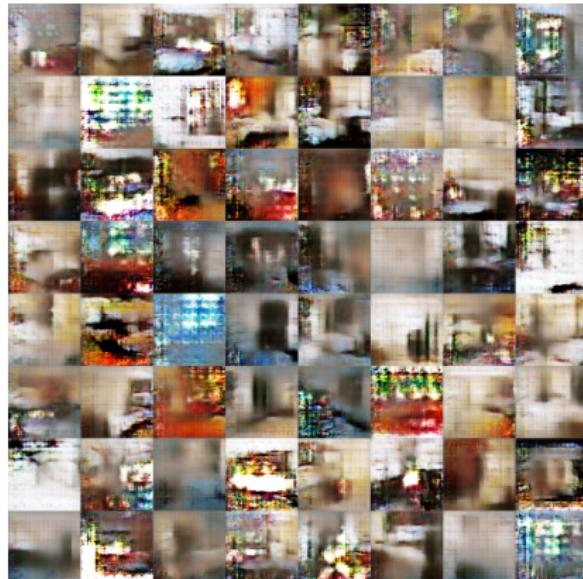
WGAN samples, $f = 64$,
FID=41, KID=4 44/53

Results: what does MMD buy you?

- Critic features from DCGAN: an f -filter critic has f , $2f$, $4f$ and $8f$ convolutional filters in layers 1-4. LSUN 64×64 .



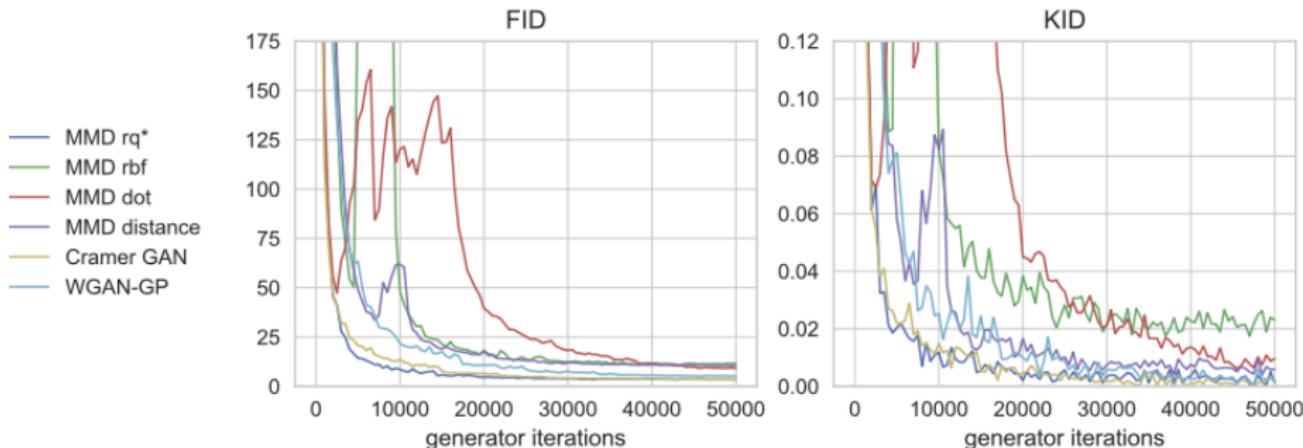
MMD GAN samples, $f = 16$,
FID=86, KID=9



WGAN samples, $f = 16$,
 $f = 64$, FID=293, KID=37
44/53

The kernel inception distance (KID)

Faster training: performance scores vs generator iterations on MNIST



Results: celebrity faces 160×160

KID (FID)
scores:

- Sobolev GAN:
14 (20)
- SN-GAN:
18 (28)
- Old MMD
GAN:
13 (21)
- SMMD GAN:
6 (12)

202 599 face images, re-sized and cropped to 160 × 160

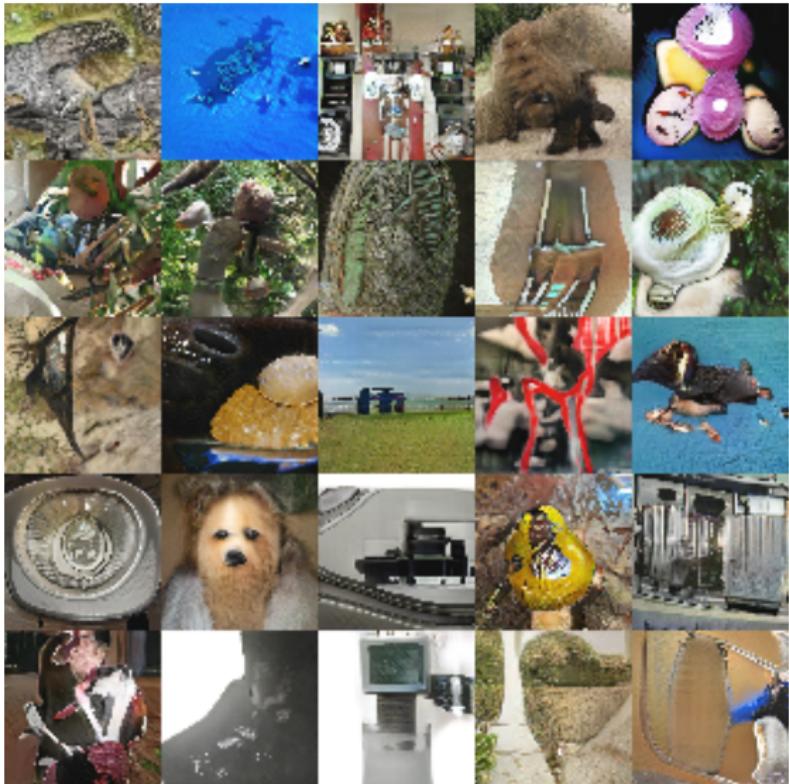


Results: imagenet 64×64

KID (FID)
scores:

- BGAN:
47 (44)
- SN-GAN:
44 (48)
- SMMD GAN:
35 (37)

ILSVRC2012 (ImageNet)
dataset, 1 281 167 im-
ages, resized to 64 × 64.
Around 20 000 classes.

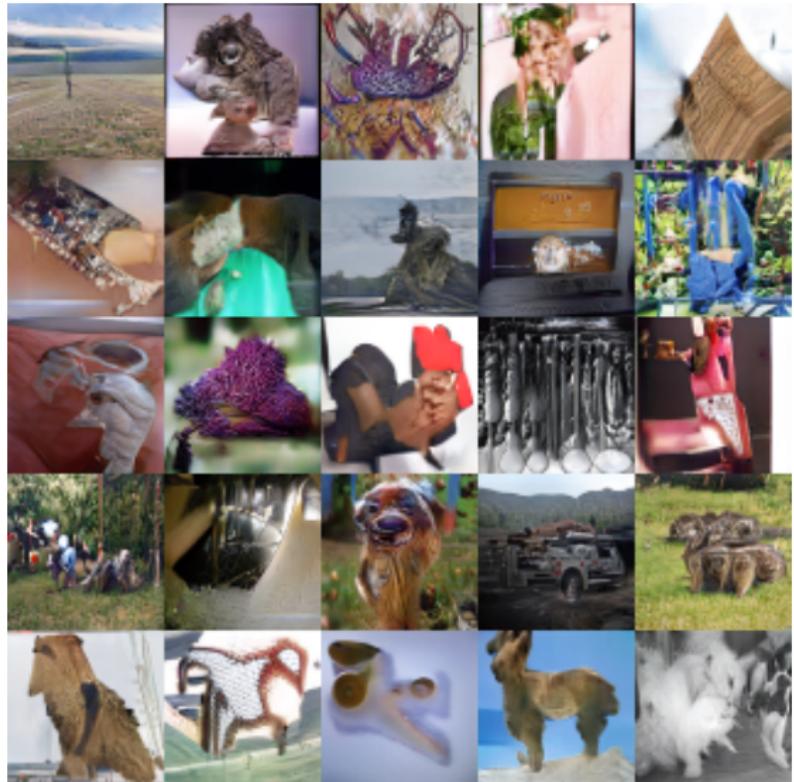


Results: imagenet 64×64

KID (FID)
scores:

- BGAN:
47 (44)
- SN-GAN:
44 (48)
- SMMD GAN:
35 (37)

ILSVRC2012 (ImageNet)
dataset, 1 281 167 im-
ages, resized to 64 × 64.
Around 20 000 classes.

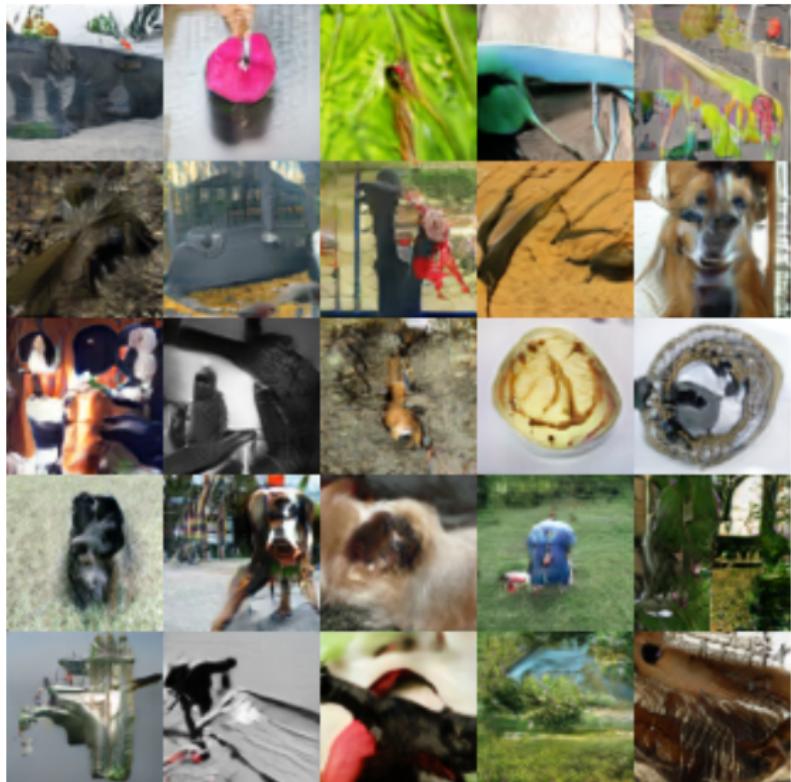


Results: imagenet 64×64

KID (FID)
scores:

- BGAN:
47 (44)
- SN-GAN:
44 (48)
- SMMD GAN:
35 (37)

ILSVRC2012 (ImageNet)
dataset, 1 281 167 im-
ages, resized to 64 × 64.
Around 20 000 classes.



Summary

- MMD critic gives state-of-the-art performance for GAN training (FID and KID)
 - use convolutional input features
 - train with gradient regulariser
- Faster training, simpler critic network
- Reasons for good performance:
 - Unlike WGAN-GP, MMD loss still a valid critic when features not optimal
 - Kernel features do some of the “work”, so simpler h_ψ features possible.
 - Better gradient/feature regulariser gives better critic

Code for “Demystifying MMD GANs,” ICLR 2018, including KID score: <https://github.com/mbinkowski/MMD-GAN>

Code for new SMMD:

<https://github.com/MichaelArbel/Scaled-MMD-GAN>

Questions?



DEMYSTIFYING MMD GANs

Mikolaj Bińkowski*
Department of Mathematics
Imperial College London
mikbinkowski@gmail.com

Dougal J. Sutherland; Michael Arbel & Arthur Gretton
Gatsby Computational Neuroscience Unit
University College London
{dougal,michael.n.arbel,arthur.gretton}@gmail.com

[arXiv.org > stat > arXiv:1805.11565](https://arxiv.org/abs/1805.11565)

Statistics > Machine Learning

On gradient regularizers for MMD GANs

Michael Arbel, Dougal J. Sutherland, [Mikolaj Bińkowski](#), Arthur Gretton

(Submitted on 29 May 2018)

Gradient bias for MMD GANs

MMD and WGAN-GP have unbiased gradients

Recall definitions:

- Generator G_ψ with parameters ψ . Samples $\textcolor{red}{Y} \sim G_\psi(\textcolor{red}{Z})$ where $\textcolor{red}{Z} \sim \textcolor{red}{R}$
- Critic features h_θ with parameters θ .

Subject to mild conditions on the mappings h_θ and G_ψ , the kernel k , and the distributions P and R , for μ -almost all ψ, θ where μ is Lebesgue,

$$\mathbb{E}_{\substack{X \sim P \\ Z \sim R}} [\partial_{\psi, \theta} k(h_\theta(X), h_\theta(G_\psi(Z)))] = \partial_{\psi, \theta} \mathbb{E}_{\substack{X \sim P \\ Z \sim R}} [k(h_\theta(X), h_\theta(G_\psi(Z)))] .$$

and thus MMD gradients unbiased.

Same true for WGAN-GP.

MMD and WGAN-GP have unbiased gradients

Recall definitions:

- Generator G_ψ with parameters ψ . Samples $\textcolor{red}{Y} \sim G_\psi(\textcolor{red}{Z})$ where $\textcolor{red}{Z} \sim \textcolor{red}{R}$
- Critic features h_θ with parameters θ .

Subject to mild conditions on the mappings h_θ and G_ψ , the kernel k , and the distributions $\textcolor{blue}{P}$ and $\textcolor{red}{R}$, for μ -almost all ψ, θ where μ is Lebesgue,

$$\mathbf{E}_{\substack{X \sim \textcolor{blue}{P} \\ Z \sim \textcolor{red}{R}}} [\partial_{\psi, \theta} k(h_\theta(X), h_\theta(G_\psi(Z)))] = \partial_{\psi, \theta} \mathbf{E}_{\substack{X \sim \textcolor{blue}{P} \\ Z \sim \textcolor{red}{R}}} [k(h_\theta(X), h_\theta(G_\psi(Z)))] .$$

and thus MMD gradients unbiased.

Same true for WGAN-GP.

Bias of MMD GAN critic

Gradient bias when critic trained on a separate dataset?

Recall definition of MMD for P vs Q

$$MMD(P, Q; F) := \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$(F = \text{unit ball in RKHS } \mathcal{F})$

Define f_{tr} as discriminator witness trained on $\{x_i^{\text{tr}}\}_{i=1}^m \stackrel{\text{i.i.d.}}{\sim} P$,
 $\{y_i^{\text{tr}}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} Q$.

Then

$$[\mathbb{E}_P f_{tr}(X) - \mathbb{E}_Q f_{tr}(Y)] \leq MMD(P, Q; F)$$

Downwards bias. Unless bias is in f_{tr} constant, biased gradients too.

Same true for WGAN-GP.

Bias of MMD GAN critic

Gradient bias when critic trained on a separate dataset?

Recall definition of MMD for P vs Q

$$MMD(P, Q; F) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(F = \text{unit ball in RKHS } \mathcal{F})$

Define f_{tr} as discriminator witness trained on $\{x_i^{\text{tr}}\}_{i=1}^m \stackrel{\text{i.i.d.}}{\sim} P$,
 $\{y_i^{\text{tr}}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} Q$.

Then

$$[\mathbf{E}_P f_{tr}(X) - \mathbf{E}_Q f_{tr}(Y)] \leq MMD(P, Q; F)$$

Downwards bias. Unless bias is in f_{tr} constant, biased gradients too.

Same true for WGAN-GP.

Bias of MMD GAN critic

Gradient bias when critic trained on a separate dataset?

Recall definition of MMD for P vs Q

$$MMD(P, Q; F) := \sup_{\|f\| \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(F = \text{unit ball in RKHS } \mathcal{F})$

Define f_{tr} as discriminator witness trained on $\{x_i^{\text{tr}}\}_{i=1}^m \stackrel{\text{i.i.d.}}{\sim} P$,
 $\{y_i^{\text{tr}}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} Q$.

Then

$$[\mathbf{E}_P f_{tr}(X) - \mathbf{E}_Q f_{tr}(Y)] \leq MMD(P, Q; F)$$

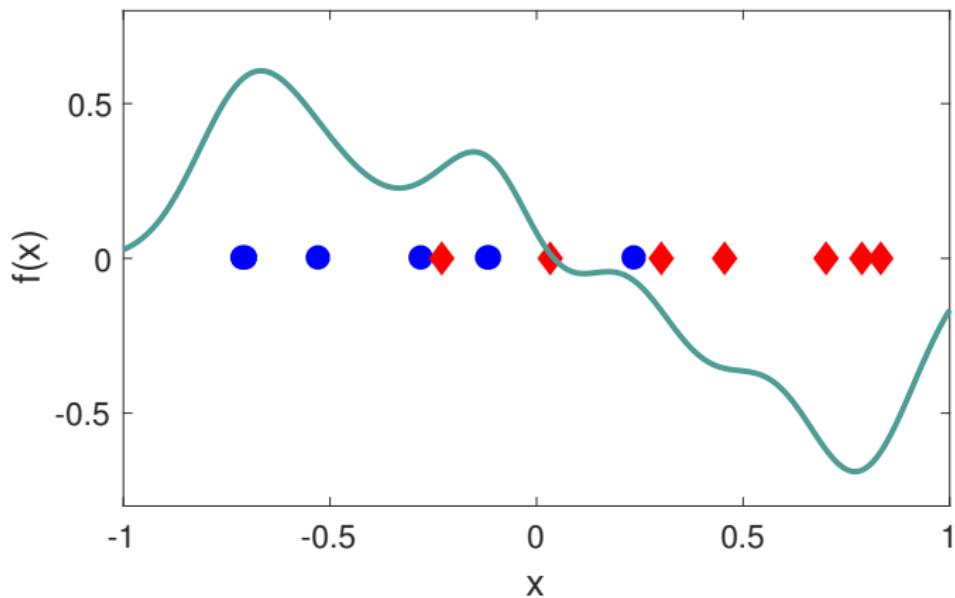
Downwards bias. Unless bias is in f_{tr} constant, biased gradients too.

Same true for WGAN-GP.

Bias of MMD GAN critic

Training minibatch critic function f_{tr}

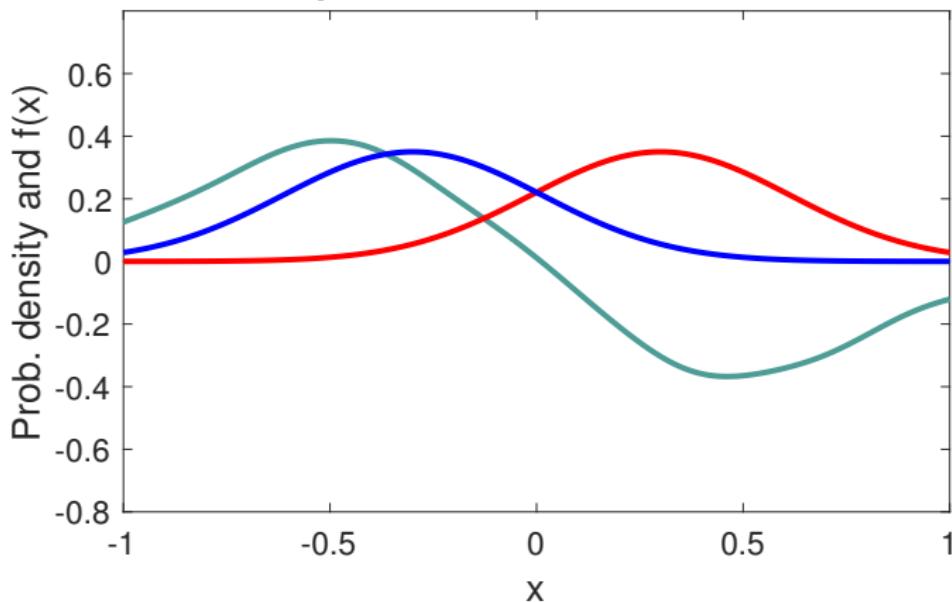
Trained witness function f_{tr}



Bias of MMD GAN critic

Population critic function f^*

Population witness function f^*



Bias of MMD GAN critic

Bias in MMD vs training minibatch size:

