

Probability Divergences and Generative Models

Arthur Gretton

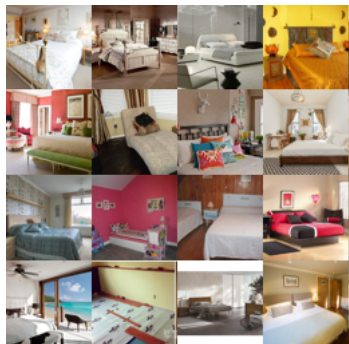


Gatsby Computational Neuroscience Unit,
University College London

PAISS 2021

Training generative models

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



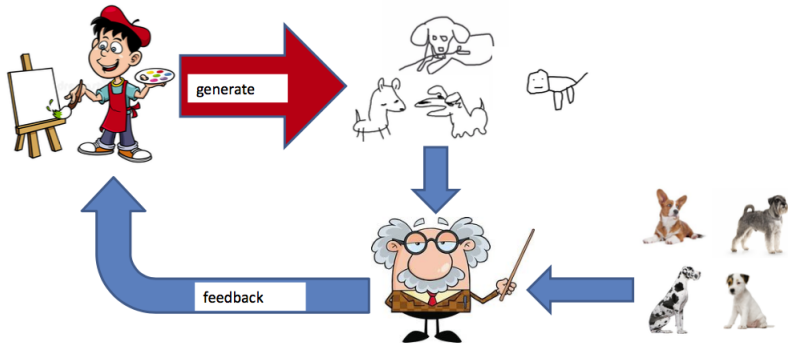
LSUN bedroom samples P



Generated Q , MMD GAN

Role of divergence $D(P, Q)$?

Reminder: generative adversarial network



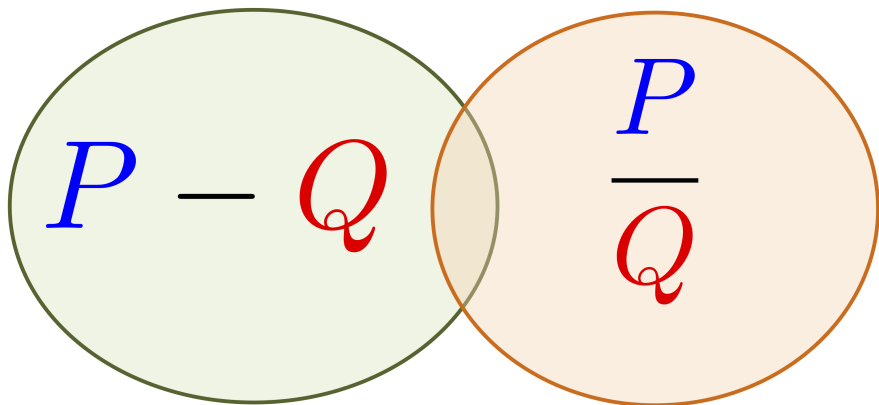
Outline

- Integral probability metrics (MMD, Wasserstein)
- ϕ -divergences (f -divergences) and a variational lower bound (KL)
- Generalized energy-based models
 - “Like a GAN” but incorporating **critic** into sample generation
 - Performs better than using **generator** alone

Arbel, Zhou, G., Generalized Energy Based Models (ICLR 2021)

Divergence measures (critics)

Divergences



Divergences

Integral prob. metrics

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

ϕ -divergences

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

The Integral Probability Metrics

Integral prob. metrics

wasserstein

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

MMD

ϕ -divergences

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

Wasserstein distance

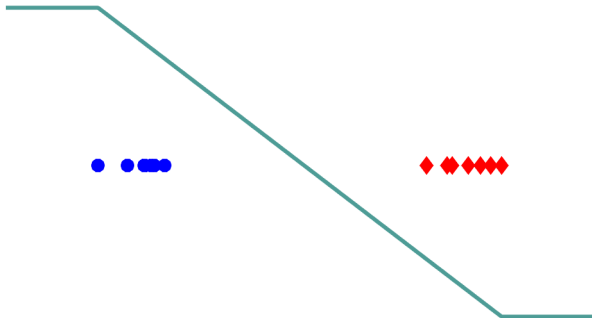


A helpful critic:

$$W_1(P, Q) = \sup_{\|f\|_L \leq 1} E_P f(X) - E_Q f(Y).$$

$$\|f\|_L := \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$$

$$W_1 = 0.88$$



Santambrogio, Optimal Transport for Applied Mathematicians (2015, Section 5.4)

G Peyré, M Cuturi, Computational Optimal Transport (2019)

M. Cuturi, J. Solomon, NeurIPS tutorial (2017)

Wasserstein distance

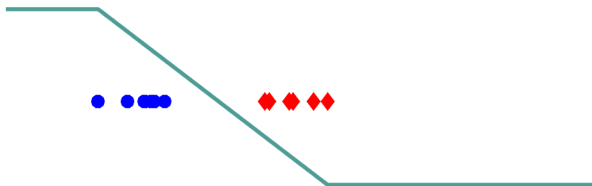


A helpful critic:

$$W_1(P, Q) = \sup_{\|f\|_L \leq 1} E_P f(X) - E_Q f(Y).$$

$$\|f\|_L := \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$$

$$W_1 = 0.65$$



Santambrogio, Optimal Transport for Applied Mathematicians (2015, Section 5.4)

G Peyré, M Cuturi, Computational Optimal Transport (2019)

M. Cuturi, J. Solomon, NeurIPS tutorial (2017)

The Maximum Mean Discrepancy

Maximum mean discrepancy: smooth function for P vs Q

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

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

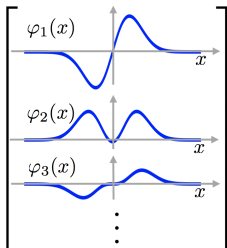
The Maximum Mean Discrepancy

Maximum mean discrepancy: smooth function for P vs Q

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

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

Functions are linear combinations of features:

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

$$\|f\|_{\mathcal{F}}^2 := \sum_{i=1}^{\infty} f_i^2 \leq 1$$

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

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

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

For positive definite k ,

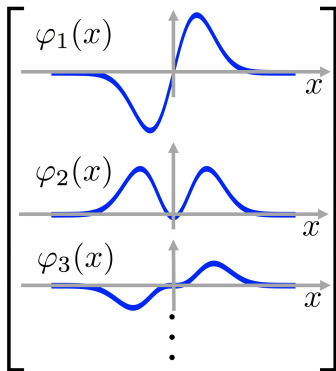
$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

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

Exponentiated quadratic kernel

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

$$\varphi(x) =$$



The MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

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

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

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

- Energy distance is a special case [Sejdinovic, Sriperumbudur, G. Fukumizu, 2013]

The MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

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

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

Expectations of functions are linear combinations of expected features

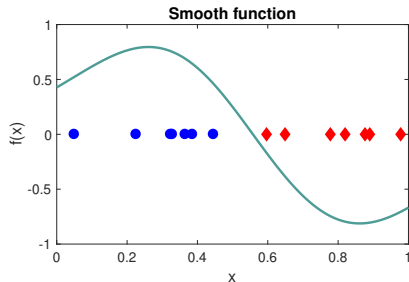
$$\mathbb{E}_P(f(X)) = \langle f, \mathbb{E}_P \varphi(X) \rangle_{\mathcal{F}} = \langle f, \mu_P \rangle_{\mathcal{F}}$$

(always true if kernel is bounded)

Integral prob. metric vs feature mean difference

The MMD:

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



Integral prob. metric vs feature mean difference

The MMD:

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

use

$$= \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

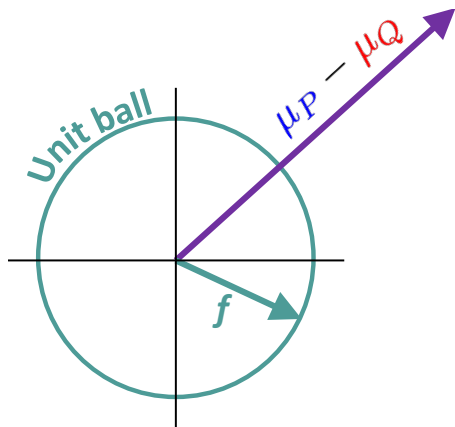
$$\mathbb{E}_P f(X) = \langle \mu_P, f \rangle_{\mathcal{F}}$$

$$= \sup_{\|f\| \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$

Integral prob. metric vs feature mean difference

The MMD:

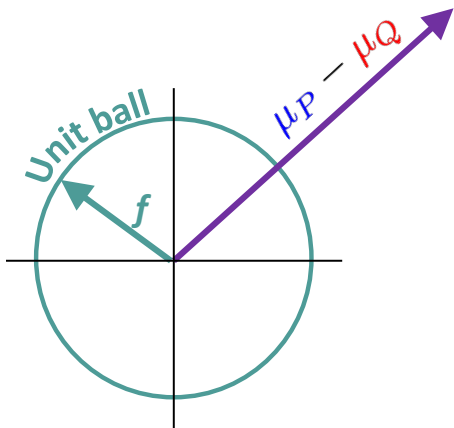
$$\begin{aligned} \text{MMD}(P, Q; F) &= \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)] \\ &= \sup_{\|f\| \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}} \end{aligned}$$



Integral prob. metric vs feature mean difference

The MMD:

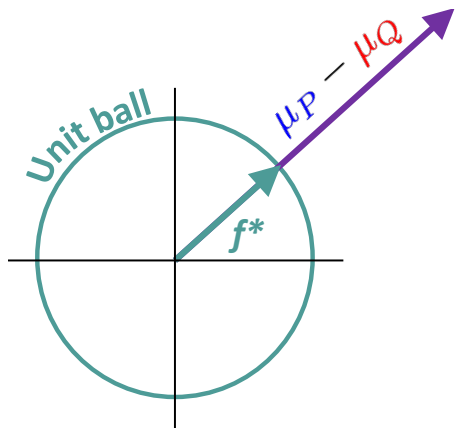
$$\begin{aligned}MMD(P, Q; F) &= \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)] \\ &= \sup_{\|f\| \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}\end{aligned}$$



Integral prob. metric vs feature mean difference

The MMD:

$$\begin{aligned} \text{MMD}(P, Q; F) &= \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)] \\ &= \sup_{\|f\| \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}} \end{aligned}$$



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

Integral prob. metric vs feature mean difference

The MMD:

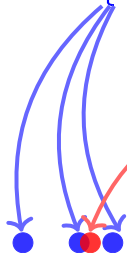
$$\begin{aligned} &MMD(P, Q; F) \\ &= \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)] \\ &= \sup_{\|f\| \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}} \\ &= \|\mu_P - \mu_Q\| \end{aligned}$$

IPM view equivalent to feature mean difference (kernel case only)

Construction of MMD witness

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

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

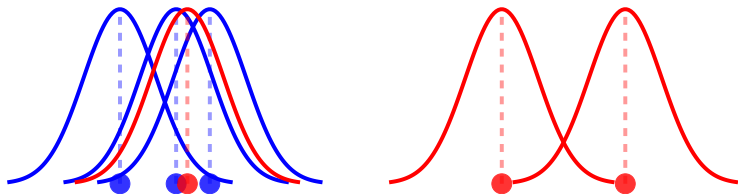


Observe $Y = \{y_1, \dots, y_n\} \sim Q$



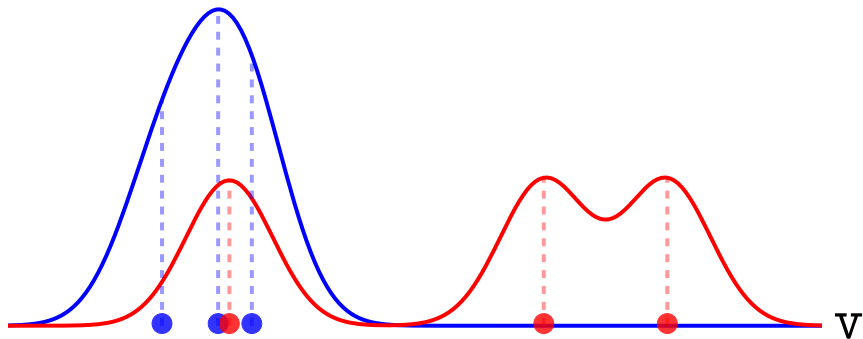
Construction of MMD witness

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



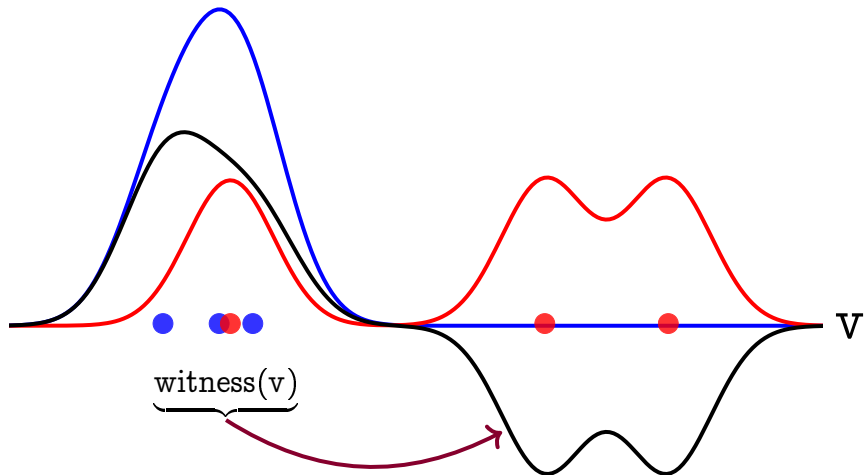
Construction of MMD witness

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

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

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

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

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

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

Don't need explicit feature coefficients $f^* := \begin{bmatrix} f_1^* & f_2^* & \dots \end{bmatrix}$

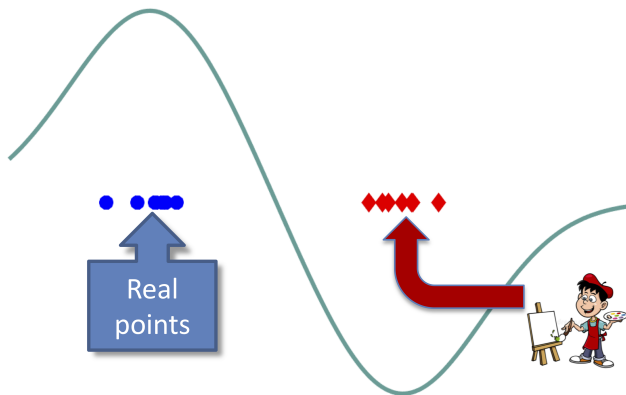
Maximum mean discrepancy



A helpful critic:

$$MMD(P, Q) = \sup_{\|f\|_{\mathcal{F}} \leq 1} E_P f(X) - E_Q f(Y).$$

MMD=1.8



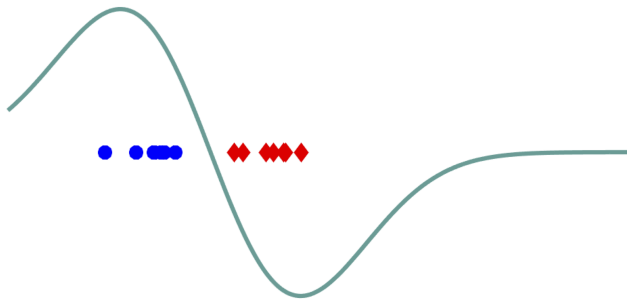
Maximum mean discrepancy



A helpful critic:

$$MMD(P, Q) = \sup_{\|f\|_{\mathcal{F}} \leq 1} E_P f(X) - E_Q f(Y)$$

MMD=1.1



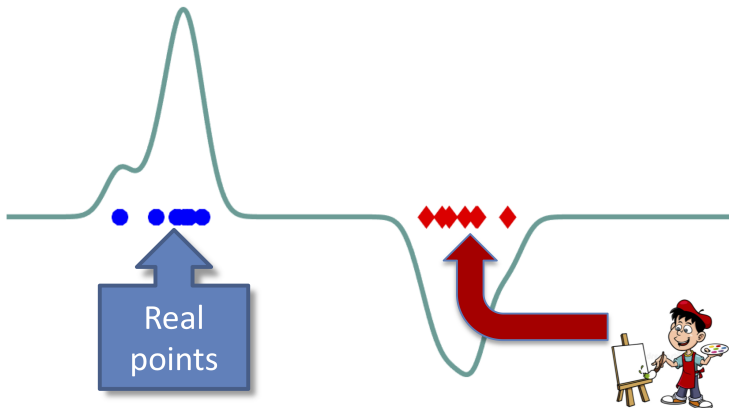
Maximum mean discrepancy



An **unhelpful** critic:

$MMD(P, Q)$ with a narrow kernel.

MMD=0.64



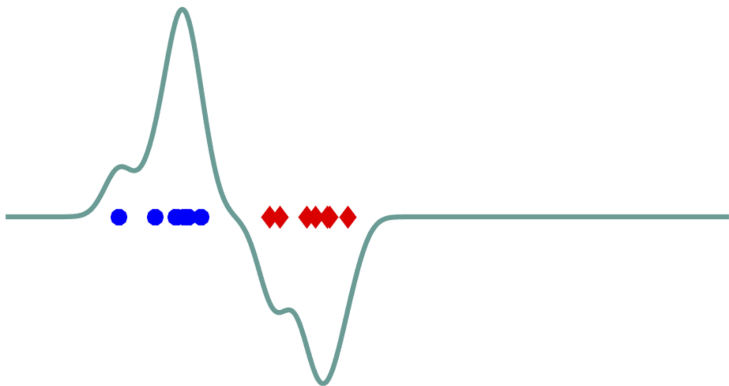
Maximum mean discrepancy



An **unhelpful** critic:

$MMD(P, Q)$ with a narrow kernel.

MMD=0.64



The ϕ -divergences

Integral prob. metrics

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

ϕ -divergences

Hellinger

KL

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

Pearson chi²

The ϕ -divergences

Define the ϕ -divergence (f -divergence):

$$D_{\phi}(P, Q) = \int \phi \left(\frac{p(z)}{q(z)} \right) q(z) dz$$

where ϕ is convex, lower-semicontinuous, $\phi(1) = 0$.

■ **Example:** $\phi(u) = u \log(u)$ gives KL divergence,

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &= \int \left(\frac{p(z)}{q(z)} \right) \log \left(\frac{p(z)}{q(z)} \right) q(z) dz \end{aligned}$$

The ϕ -divergences

Define the ϕ -divergence (f -divergence):

$$D_{\phi}(P, Q) = \int \phi \left(\frac{p(z)}{q(z)} \right) q(z) dz$$

where ϕ is convex, lower-semicontinuous, $\phi(1) = 0$.

■ **Example:** $\phi(u) = u \log(u)$ gives KL divergence,

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &= \int \left(\frac{p(z)}{q(z)} \right) \log \left(\frac{p(z)}{q(z)} \right) q(z) dz \end{aligned}$$

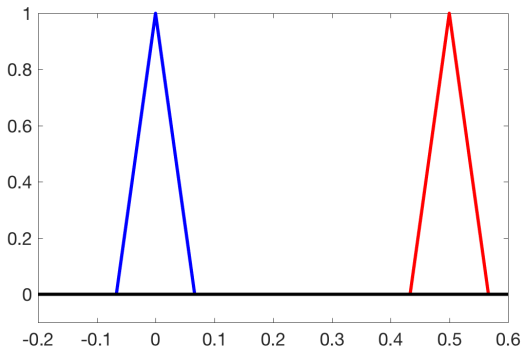
Are ϕ -divergences good critics?



Simple example: disjoint support.

Goodfellow et al. (NeurIPS 2014), Arjovsky and Bottou [ICLR 2017]

$$D_{KL}(P, Q) = \infty \quad D_{JS}(P, Q) = \log 2$$



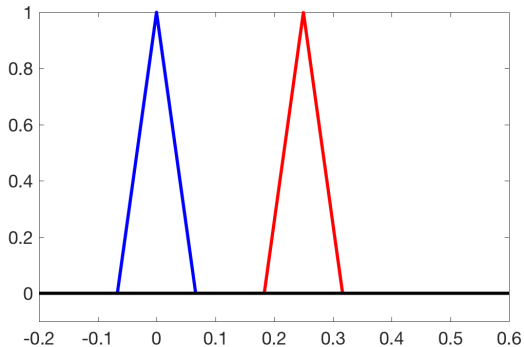
Are ϕ -divergences good critics?



Simple example: disjoint support.

Goodfellow et al. (NeurIPS 2014), Arjovsky and Bottou [ICLR 2017]

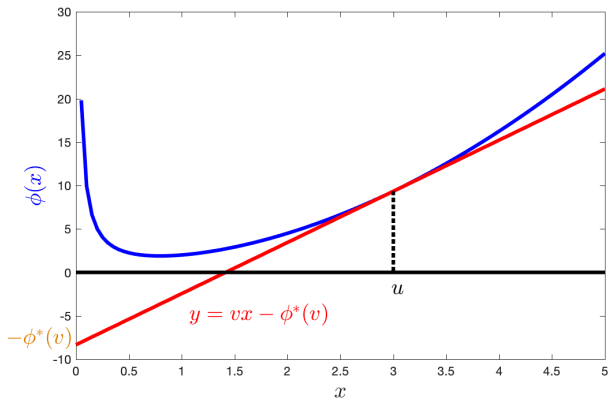
$$D_{KL}(P, Q) = \infty \quad D_{JS}(P, Q) = \log 2$$



ϕ -divergences in practice

Background: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$



- $\phi^*(v)$ is negative intercept of tangent to ϕ with slope v

ϕ -divergences in practice

Background: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$

■ For a convex l.s.c. ϕ we have

$$\phi^{**}(x) = \phi(x) = \sup_{v \in \mathbb{R}} \{xv - \phi^*(v)\}$$

ϕ -divergences in practice

Background: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$

■ For a convex l.s.c. ϕ we have

$$\phi^{**}(x) = \phi(x) = \sup_{v \in \mathbb{R}} \{xv - \phi^*(v)\}$$

■ **KL divergence:**

$$\phi(x) = x \log(x) \quad \phi^*(v) = \exp(v - 1)$$

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$D_{\phi}(P, Q) = \int q(z) \phi\left(\frac{p(z)}{q(z)}\right) dz$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned} D_\phi(P, Q) &= \int q(z) \phi\left(\frac{p(z)}{q(z)}\right) dz \\ &= \int q(z) \underbrace{\sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z) \right)}_{\phi\left(\frac{p(z)}{q(z)}\right)} dz \end{aligned}$$

$\phi^*(v)$ is dual of $\phi(v)$.

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned}D_{\phi}(P, Q) &= \int q(z) \phi\left(\frac{p(z)}{q(z)}\right) dz \\&= \int q(z) \sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z)\right) \\&\geq \sup_{f \in \mathcal{H}} \mathbb{E}_P f(X) - \mathbb{E}_Q \phi^*(f(Y))\end{aligned}$$

(restrict the function class)

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned}D_{\phi}(P, Q) &= \int q(z) \phi\left(\frac{p(z)}{q(z)}\right) dz \\&= \int q(z) \sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z)\right) \\&\geq \sup_{f \in \mathcal{H}} \mathbb{E}_P f(X) - \mathbb{E}_Q \phi^*(f(Y))\end{aligned}$$

(restrict the function class)

Bound tight when:

$$f^{\diamond}(z) = \partial \phi\left(\frac{p(z)}{q(z)}\right)$$

if ratio defined.

Case of the KL

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

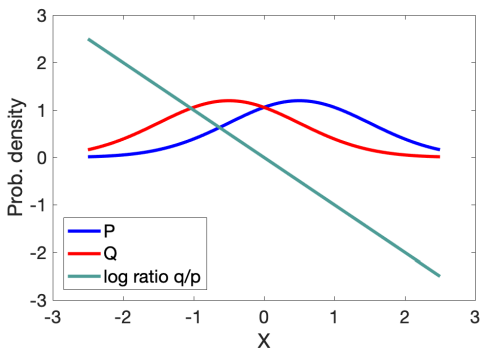
Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \underbrace{\mathbb{E}_Q \exp(-f(Y))}_{\phi^*(-f(Y)+1)} \end{aligned}$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$
$$\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \mathbb{E}_Q \exp(-f(Y))$$



Bound tight when:

$$f^\diamond(z) = -\log \frac{p(z)}{q(z)}$$

if ratio defined.

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

$$\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \mathbb{E}_Q \exp(-f(Y))$$

$$\approx \sup_{f \in \mathcal{H}} \left[-\frac{1}{n} \sum_{j=1}^n f(x_j) - \frac{1}{n} \sum_{i=1}^n \exp(-f(y_i)) \right] + 1$$

$x_i \stackrel{\text{i.i.d.}}{\sim} P$

$y_i \stackrel{\text{i.i.d.}}{\sim} Q$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \mathbb{E}_Q \exp(-f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[-\frac{1}{n} \sum_{j=1}^n f(x_j) - \frac{1}{n} \sum_{i=1}^n \exp(-f(y_i)) \right] + 1 \end{aligned}$$

This is a

KL

Approximate

Lower-bound

Estimator.

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \mathbb{E}_Q \exp(-f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[-\frac{1}{n} \sum_{j=1}^n f(x_j) - \frac{1}{n} \sum_{i=1}^n \exp(-f(y_i)) \right] + 1 \end{aligned}$$

This is a

K
A
L
E

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} -\mathbb{E}_P f(X) + 1 - \mathbb{E}_Q \exp(-f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[-\frac{1}{n} \sum_{j=1}^n f(x_j) - \frac{1}{n} \sum_{i=1}^n \exp(-f(y_i)) \right] + 1 \end{aligned}$$

The KALE divergence

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Topological properties of KALE (1)

Key requirements on \mathcal{H} and \mathcal{X} :

- Compact domain \mathcal{X} ,
- \mathcal{H} dense in the space $C(\mathcal{X})$ of continuous functions on \mathcal{X} wrt $\|\cdot\|_\infty$.
- If $f \in \mathcal{H}$ then $-f \in \mathcal{H}$ and $cf \in \mathcal{H}$ for $0 \leq c \leq C_{\max}$.

Theorem: $KALE(P, Q; \mathcal{H}) \geq 0$ and $KALE(P, Q; \mathcal{H}) = 0$ iff $P = Q$.

Zhang, Liu, Zhou, Xu, and He. "On the Discrimination-Generalization Tradeoff in GANs"
(ICLR 2018, Corollary 2.4; Theorem B.1)
Arbel, Liang, G. (ICLR 2021, Proposition 1)

Topological properties of KALE (1)

Key requirements on \mathcal{H} and \mathcal{X} :

- Compact domain \mathcal{X} ,
- \mathcal{H} dense in the space $C(\mathcal{X})$ of continuous functions on \mathcal{X} wrt $\|\cdot\|_\infty$.
- If $f \in \mathcal{H}$ then $-f \in \mathcal{H}$ and $cf \in \mathcal{H}$ for $0 \leq c \leq C_{\max}$.

Theorem: $KALE(P, Q; \mathcal{H}) \geq 0$ and $KALE(P, Q; \mathcal{H}) = 0$ iff $P = Q$.

\mathcal{H} dense in $C(\mathcal{X})$ for $\mathcal{X} \subset \mathbb{R}^d$ when:

$$\mathcal{H} = \text{span}\{\sigma(w^\top x + b) : [w, b] \in \Theta\}$$

$$\sigma(u) = \max\{u, 0\}^\alpha, \alpha \in \mathbb{N}, \text{ and } \{\lambda\theta : \lambda \geq 0, \theta \in \Theta\} = \mathbb{R}^{d+1}.$$

Zhang, Liu, Zhou, Xu, and He. "On the Discrimination-Generalization Tradeoff in GANs"

(ICLR 2018, Corollary 2.4; Theorem B.1)

Arbel, Liang, G. (ICLR 2021, Proposition 1)

Topological properties of KALE (2)

Additional requirement: all functions in \mathcal{H} Lipschitz in their inputs with constant L

Theorem: $KALE(P, Q^n; \mathcal{H}) \rightarrow 0$ iff $Q^n \rightarrow P$ under the weak topology.

Topological properties of KALE (2)

Additional requirement: all functions in \mathcal{H} Lipschitz in their inputs with constant L

Theorem: $KALE(P, Q^n; \mathcal{H}) \rightarrow 0$ iff $Q^n \rightarrow P$ under the weak topology.

Partial proof idea:

$$\begin{aligned} KALE(P, Q; \mathcal{H}) &= - \int f dP - \int \exp(-f) dQ + 1 \\ &= \int f(x) dQ(x) - \int f(x') dP(x') \\ &\quad - \underbrace{\int (\exp(-f) + f - 1) dQ}_{\geq 0} \\ &\leq \int f(x) dQ(x) - \int f(x') dP(x') \leq LW_1(P, Q) \end{aligned}$$

Empirical properties of KALE



$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} -E_P f(X) - E_Q \exp(-f(Y)) + 1$$

$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$

$$\|w\|_{\mathcal{H}}^2 \text{ penalized :}$$

Empirical properties of KALE



$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} -E_P f(X) - E_Q \exp(-f(Y)) + 1$$

$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$

$$\|w\|_{\mathcal{H}}^2 \text{ penalized : KALE smoothie}$$

Empirical properties of KALE

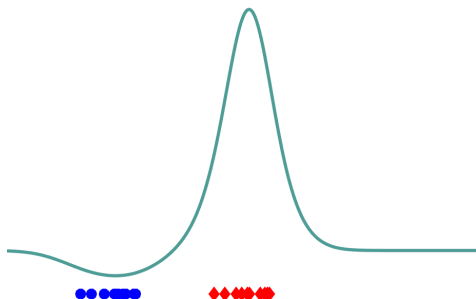


$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} -E_P f(X) - E_Q \exp(-f(Y)) + 1$$

$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$

$\|w\|_{\mathcal{H}}^2$ penalized : KALE smoothie

$$KALE(Q, P; \mathcal{H}) = 0.18$$



Empirical properties of KALE

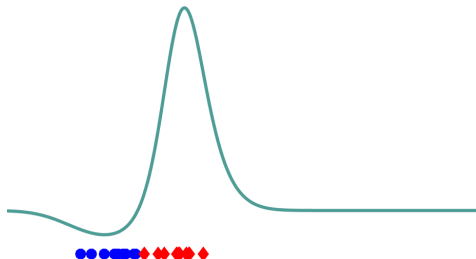


$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} -E_P f(X) - E_Q \exp(-f(Y)) + 1$$

$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$

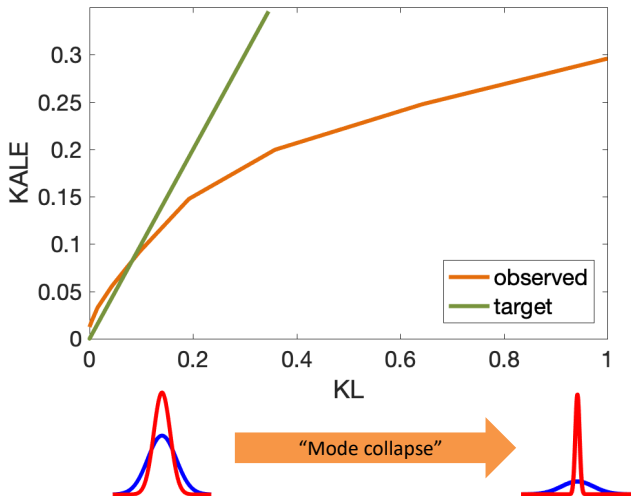
$\|w\|_{\mathcal{H}}^2$ penalized : KALE smoothie

$$KALE(Q, P; \mathcal{H}) = 0.12$$



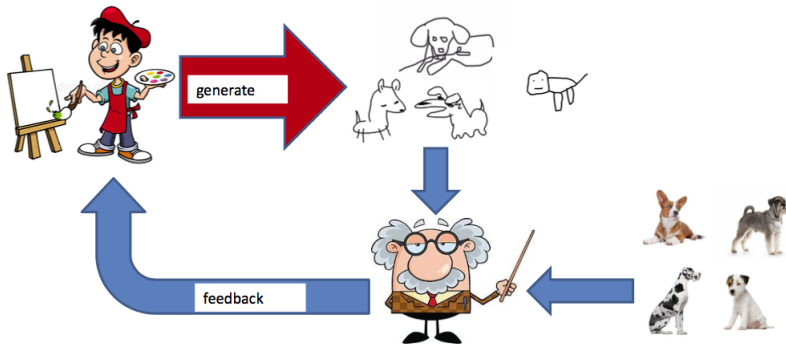
The KALE smoothie and “mode collapse”

- Two Gaussians with same means, different variance

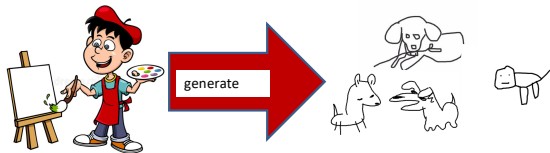


Generalized Energy-Based Models

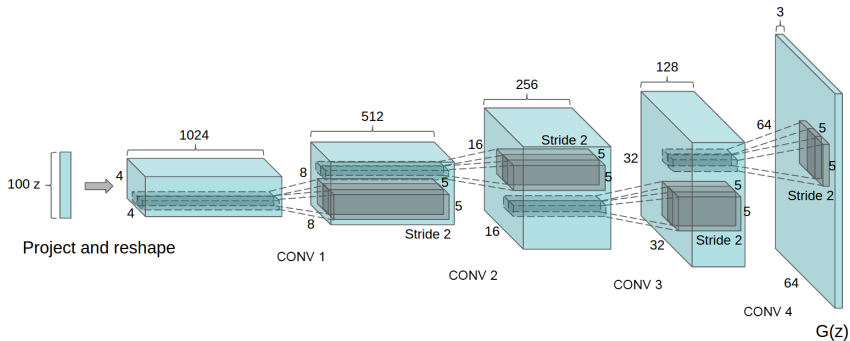
Visual notation: GAN setting



Visual notation: GAN setting



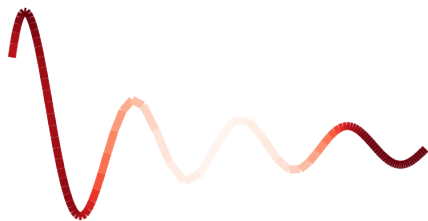
Reminder: the generator



Radford, Metz, Chintala, ICLR 2016

Generalized energy-based models: illustration

Target distribution P



$$z \sim \text{Unif}[0, 1]$$

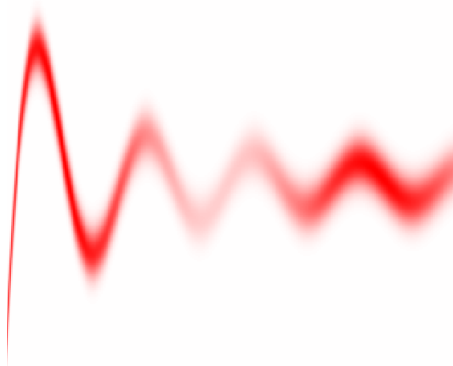
$$\tilde{z} = \tau(z)$$

$$X = G_{\theta^*}(\tilde{z}), \quad X_1 = \tilde{z}$$

Example thanks to M. Arbel

Generalized energy-based models: illustration

EBM approximation to target:

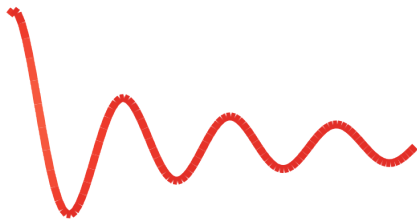


$$p(X) \propto \exp(-E(X))$$
$$E(X) = \frac{1}{2\sigma^2} \|G_\theta(X_1) - X\|^2 + A_\theta(X_1)$$

Example thanks to M. Arbel

Generalized energy-based models: illustration

GAN (generator) distribution Q_θ



Generator

$$z \sim \text{unif}[0, 1]$$

$$X = B_\theta(z)$$

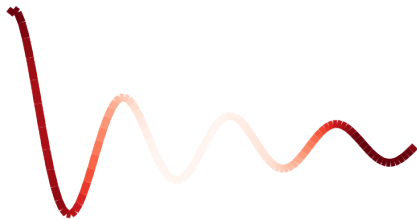
Critic

$$MLP(X)$$

Example thanks to M. Arbel

Generalized energy-based models: illustration

Mass of GEBM corrected by critic



Generator

$$z \sim \text{unif}[0, 1]$$

$$X = B_{\theta}(z)$$

Re-weight using importance weights defined by energy:

$$w(x) \propto \exp(-E(x))$$

Example thanks to M. Arbel

Generalized energy-based models

Define a model $Q_{B_\theta, E}$ as follows:

- Sample from **generator** with parameters θ

$$X \sim Q_\theta \iff X = B_\theta(Z), \quad Z \sim \eta$$

- Reweight the samples according to importance weights:

$$f_{Q, E}(x) = \frac{\exp(-E(x))}{Z_{Q, E}}, \quad Z_{Q, E} = \int \exp(-E(x)) dQ_\theta(x),$$

where $E \in \mathcal{E}$, the energy function class.

$f_{Q, E}(x)$ is Radon-Nikodym derivative of $Q_{B_\theta, E}$ wrt Q_θ .

- When Q_θ has density wrt Lebesgue on \mathcal{X} , this is a standard energy-based model.

Fitting GEBMs

Fit the model using **Generalized Log-Likelihood**:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

- When $KL(P, Q_\theta)$ well defined, above is **Donsker-Varadhan** lower bound on KL
 - tight when $E(z) = -\log(p(z)/q(z))$.
- However, **Generalized Log-Likelihood** still defined when P and Q_θ mutually singular!

KALE and the energy function

Fit the model using Generalized Log-Likelihood:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

KALE and the energy function

Fit the model using Generalized Log-Likelihood:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

From convexity of exponential,

$$- \log \int \exp(-E) dQ_\theta \geq -c - e^{-c} \int \exp(-E) dQ_\theta + 1$$

tight whenever $c = \log \int \exp(-E) dQ_\theta$.

KALE and the energy function

Fit the model using Generalized Log-Likelihood:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

From convexity of exponential,

$$- \log \int \exp(-E) dQ_\theta \geq -c - e^{-c} \int \exp(-E) dQ_\theta + 1$$

tight whenever $c = \log \int \exp(-E) dQ_\theta$.

Generalized Log-Likelihood has the lower bound:

$$\begin{aligned} \mathcal{L}_{P,Q}(E) &\geq - \int (E + c) dP - \int \exp(-E - c) dQ_\theta + 1 \\ &:= \mathcal{F}(P, Q_\theta; \mathcal{E} + \mathbb{R}) \end{aligned}$$

KALE and the energy function

Fit the model using Generalized Log-Likelihood:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

From convexity of exponential,

$$- \log \int \exp(-E) dQ_\theta \geq -c - e^{-c} \int \exp(-E) dQ_\theta + 1$$

tight whenever $c = \log \int \exp(-E) dQ_\theta$.

Generalized Log-Likelihood has the lower bound:

$$\begin{aligned} \mathcal{L}_{P,Q}(E) &\geq - \int (E + c) dP - \int \exp(-E - c) dQ_\theta + 1 \\ &:= \mathcal{F}(P, Q_\theta; \mathcal{E} + \mathbb{R}) \end{aligned}$$

This is the KALE with function class $\mathcal{E} + \mathbb{R}$.

KALE and the energy function

Fit the model using Generalized Log-Likelihood:

$$\mathcal{L}_{P,Q}(E) := \int \log(f_{Q,E}) dP = - \int E dP - \log \int \exp(-E) dQ_\theta$$

From convexity of exponential,

$$- \log \int \exp(-E) dQ_\theta \geq -c - e^{-c} \int \exp(-E) dQ_\theta + 1$$

tight whenever $c = \log \int \exp(-E) dQ_\theta$.

Generalized Log-Likelihood has the lower bound:

$$\begin{aligned} \mathcal{L}_{P,Q}(E) &\geq - \int (E + c) dP - \int \exp(-E - c) dQ_\theta + 1 \\ &:= \mathcal{F}(P, Q_\theta; \mathcal{E} + \mathbb{R}) \end{aligned}$$

Jointly maximizing yields the maximum likelihood energy E^* and corresponding $c^* = \log \int \exp(-E) dQ_\theta$.

Training the base measure (generator)

Recall the generator:

$$X = B_{\theta}(Z), \quad Z \sim \eta$$

Define: $\mathcal{K}(\theta) := \mathcal{F}(P, Q_{\theta}; \mathcal{E} + \mathbb{R})$

Training the base measure (generator)

Recall the **generator**:

$$X = B_{\theta}(Z), \quad Z \sim \eta$$

Define: $\mathcal{K}(\theta) := \mathcal{F}(P, Q_{\theta}; \mathcal{E} + \mathbb{R})$

Theorem: \mathcal{K} is lipschitz and differentiable for almost all $\theta \in \Theta$ with:

$$\nabla \mathcal{K}(\theta) = Z_{Q, E^*}^{-1} \int \nabla_x E^*(B_{\theta}(z)) \nabla_{\theta} B_{\theta}(z) \exp(-E^*(B_{\theta}(z))) \eta(z) dz.$$

where E^* achieves supremum in $\mathcal{F}(P, Q; \mathcal{E} + \mathbb{R})$.

Training the base measure (generator)

Recall the **generator**:

$$X = B_{\theta}(Z), \quad Z \sim \eta$$

Define: $\mathcal{K}(\theta) := \mathcal{F}(P, Q_{\theta}; \mathcal{E} + \mathbb{R})$

Theorem: \mathcal{K} is lipschitz and differentiable for almost all $\theta \in \Theta$ with:

$$\nabla \mathcal{K}(\theta) = Z_{Q, E^*}^{-1} \int \nabla_x E^*(B_{\theta}(z)) \nabla_{\theta} B_{\theta}(z) \exp(-E^*(B_{\theta}(z))) \eta(z) dz.$$

where E^* achieves supremum in $\mathcal{F}(P, Q; \mathcal{E} + \mathbb{R})$.

Assumptions:

- Functions in \mathcal{E} parametrized by $\psi \in \Psi$, where Ψ compact,
 - jointly continuous w.r.t. (ψ, x) , L -lipschitz and L -smooth w.r.t. x .
- $(\theta, z) \mapsto B_{\theta}(z)$ jointly continuous wrt (θ, z) , $z \mapsto B_{\theta}(z)$ uniformly Lipschitz w.r.t. z , lipschitz and smooth wrt θ (see paper: constants depend on z)

Sampling from the model

Consider end-to-end model $Q_{B_\theta, E}$, where recall that $X = B_\theta(Z)$, $Z \sim \eta$,

$$f_{B, E}(x) := \frac{\exp(-E(x))}{Z_{Q, E}}$$

Sampling from the model

Consider end-to-end model $Q_{B_\theta, E}$, where recall that $X = B_\theta(Z)$, $Z \sim \eta$,

$$f_{B, E}(x) := \frac{\exp(-E(x))}{Z_{Q, E}}$$

For a test function g ,

$$\int g(x) dQ_{B, E}(x) = \int g(B(z)) f_{B, E}(B(z)) \eta(z) dz$$

Posterior latent distribution therefore

$$\nu_{B, E}(z) = \eta(z) f_{B, E}(B(z))$$

Sampling from the model

Consider end-to-end model $Q_{B_\theta, E}$, where recall that $X = B_\theta(Z)$, $Z \sim \eta$,

$$f_{B, E}(x) := \frac{\exp(-E(x))}{Z_{Q, E}}$$

For a test function g ,

$$\int g(x) dQ_{B, E}(x) = \int g(B(z)) f_{B, E}(B(z)) \eta(z) dz$$

Posterior latent distribution therefore

$$\nu_{B, E}(z) = \eta(z) f_{B, E}(B(z))$$

Sample $z \sim \nu_{B, E}$ via Langevin diffusion-derived algorithms (MALA, ULA, HMC,...) to exploit gradient information.

Generate new samples in \mathcal{X} via

$$X \sim Q_{B, E} \iff Z \sim \nu_{B, E}, \quad X = B_\theta(Z).$$

Experiments

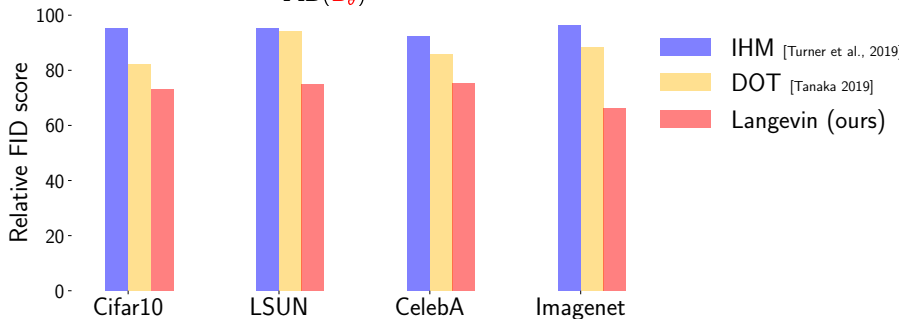
Examples: sampling at modes

Tempered GEBM Cifar10 samples at different stages of sampling using a Kinetic Langevin Algorithm (KLA). Early samples \rightarrow late samples.
Model run at low temperature ($\beta = 100$) for better quality samples.



Sampling at modes: results

The relative FID score: $\frac{\text{FID}(Q_{B_\theta, E})}{\text{FID}(B_\theta)}$



For a given generator B_θ and energy E , samples **always better** (FID score) than generator alone.

Examples: moving between modes

Tempered GEBM Cifar10 samples at different stages of sampling using KLA. Early samples \rightarrow late samples.

Model run at lower friction (but still low temperature, $\beta = 100$) for mode exploration.



Summary

- Generalized energy based model: ICLR 2021
 - End-to-end model incorporating generator and critic
 - Always better samples than generator alone.

arXiv.org > stat > arXiv:2003.05033

Statistics > Machine Learning

[Submitted on 10 Mar 2020 (v1), last revised 24 Jun 2020 (this version, v3)]

Generalized Energy Based Models

Michael Arbel, Liang Zhou, Arthur Gretton

arXiv.org > cs > arXiv:2003.06060

Search...

Help | Advan

Computer Science > Machine Learning

[Submitted on 12 Mar 2020 (v1), last revised 24 Mar 2020 (this version, v2)]

Your GAN is Secretly an Energy-based Model and You Should use Discriminator Driven Latent Sampling

Tong Che, Ruixiang Zhang, Jascha Sohl-Dickstein, Hugo Larochelle, Liam Paull, Yuan Cao, Yoshua Bengio

<https://github.com/MichaelArbel/GeneralizedEBM>

Questions?



Post-credit scene: MMD flow

From NeurIPS 2019:

Maximum Mean Discrepancy Gradient Flow

Michael Arbel

Gatsby Computational Neuroscience Unit
University College London
michael.n.arbel@gmail.com

Anna Korba

Gatsby Computational Neuroscience Unit
University College London
a.korba@ucl.ac.uk

Adil Salim

Visual Computing Center
KAUST
adil.salim@kaust.edu.sa

Arthur Gretton

Gatsby Computational Neuroscience Unit
University College London
arthur.gretton@gmail.com

Sanity check: reduction to EBM case

