

Learning with stochastic orders

Carles Domingo-Enrich^a, Yair Schiff^b and Youssef Mroueh^c

^aNew York University, ^bCornell University, ^c IBM Research AI

ICLR May, 2023

Generative adversarial networks

The goal of generative modeling is to be able to generate artificial samples from a distribution given a sample $(X_i)_{i=1}^n$ from it.

Generative adversarial networks (GANs) (Goodfellow et al., 2014) are a popular generative modeling technique where two deep neural networks, the generator g and the discriminator f , are trained adversarially.

Generative adversarial networks

The goal of generative modeling is to be able to generate artificial samples from a distribution given a sample $(X_i)_{i=1}^n$ from it.

Generative adversarial networks (GANs) (Goodfellow et al., 2014) are a popular generative modeling technique where two deep neural networks, the generator g and the discriminator f , are trained adversarially.

A common choice for the training loss (Arjovsky et al., 2017) is:

$$\min_{g \in \mathcal{G}} \max_{f \in \mathcal{F}} \mathbb{E}_{X \sim p_n} [f(X)] - \mathbb{E}_{Y \sim p_n} [f(g(Y))] ; \quad \text{where } n = \frac{1}{n} \sum_{i=1}^n x_i; \quad (1)$$

A usual failure mode of GANs is mode collapse: the generator fails to capture entire modes of the data distribution.



Figure 1: (Left) Samples from the MNIST dataset. (Right) GAN-generated samples suffering from mode collapse.

Generative adversarial networks

The goal of generative modeling is to be able to generate artificial samples from a distribution given a sample $(X_i)_{i=1}^n$ from it.

Generative adversarial networks (GANs) (Goodfellow et al., 2014) are a popular generative modeling technique where two deep neural networks, the generator g and the discriminator f , are trained adversarially.

A common choice for the training loss (Arjovsky et al., 2017) is:

$$\min_{g \in \mathcal{G}} \max_{f \in \mathcal{F}} \mathbb{E}_{X \sim p_n[f(X)]} - \mathbb{E}_{Y \sim p_n[f(g(Y))]} g; \quad \text{where } n = \frac{1}{n} \sum_{i=1}^n x_i; \quad (1)$$

A usual failure mode of GANs is mode collapse: the generator fails to capture entire modes of the data distribution.



Figure 1: (Left) Samples from the MNIST dataset. (Right) GAN-generated samples suffering from mode collapse.

Question: How can we modify the GAN objective to prevent mode collapse?
Let's look at stochastic orders first!

Can we compare probability measures beyond equality? \Rightarrow *stochastic orders*

Can we compare probability measures beyond equality? \Rightarrow *stochastic orders*

Definition (Convex or Choquet order, Ekeland and Schachermayer (2014))

For μ, ν ; μ, ν probability measures, we say that μ dominates ν in the convex order, or $\mu \succeq_{cx} \nu$, if for any convex function $u: \mathbb{R}^d \rightarrow \mathbb{R}$, we have

$$\mathbb{E}_\mu[u(x)] \geq \mathbb{E}_\nu[u(x)]:$$

Can we compare probability measures beyond equality? \Rightarrow *stochastic orders*

Definition (Convex or Choquet order, Ekeland and Schachermayer (2014))

For $\mu, \nu \in \mathcal{P}_+(\mathbb{R}^d)$ probability measures, we say that μ dominates ν in the convex order, or $\mu \succeq_{cx} \nu$, if for any convex function $u: \mathbb{R}^d \rightarrow \mathbb{R}$, we have

$$\int \mathbb{E}_\mu u(x) \geq \int \mathbb{E}_\nu u(x):$$

\succeq_{cx} is a partial order, meaning that reflexivity, antisymmetry and transitivity hold.

The space of convex functions is not the only choice to define orders (other *cones* can be considered).

Can we compare probability measures beyond equality? \Rightarrow *stochastic orders*

Definition (Convex or Choquet order, Ekeland and Schachermayer (2014))

For μ, ν ; $\mu \succeq_{cx} \nu$ probability measures, we say that μ dominates ν in the convex order, or $\mu \succeq_{cx} \nu$, if for any convex function $u : \mathbb{R}^d \rightarrow \mathbb{R}$, we have

$$\int u(x) d\mu \leq \int u(x) d\nu:$$

\succeq_{cx} is a partial order, meaning that reflexivity, antisymmetry and transitivity hold.

The space of convex functions is not the only choice to define orders (other *cones* can be considered).

The convex order in one dimension admits a characterization in terms of the integral of the CDF.

Proposition (Ekeland and Schachermayer (2014))

We have $\mu \leq_{cx} \nu$ if and only if there exists a *martingale Markov kernel* R (i.e. $\int_{\mathbb{R}^d} y dR_x(y) = x; \int_{\mathbb{R}^d} R_x d\nu = \mu$).

Proposition (Ekeland and Schachermayer (2014))

We have $\mu \preceq_{cx} \nu$ if and only if there exists a *martingale Markov kernel* R (i.e. $\int_{\mathbb{R}^d} y dR_x(y) = x; \forall x$) such that $\nu = \int_{\mathbb{R}^d} R_x d\mu$.

This characterization is difficult to check, especially in high dimensions.

Intuitively, this means that ν is more *spread out* than μ .

Variational Dominance Criterion (VDC)

Definition (Variational Dominance Criterion (VDC))

Given a bounded open convex subset $\Omega \subset \mathbb{R}^d$, a pair of Borel probability measures $\mu, \nu \in \mathcal{P}(\Omega)$, and a compact set $K \subset \mathbb{R}^d$ ($0 \in K$), define:

$$\text{VDC}_A(\mu, \nu) = \sup_{u \in A} \int_{\Omega} u d(\mu - \nu):$$

where $A = \{u : \Omega \rightarrow \mathbb{R}; u \text{ convex and } u \in K \text{ almost everywhere}\}.$

Remark that since $0 \in K$, $\text{VDC}_A(\mu, \nu) \geq 0$ for all μ, ν because the zero function belongs to the set A .

Variational Dominance Criterion (VDC)

Definition (Variational Dominance Criterion (VDC))

Given a bounded open convex subset R^d , a pair of Borel probability measures $\mu, \nu \in \mathcal{P}(R^d)$, and a compact set $K \subseteq R^d$ ($\emptyset \neq K$), define:

$$VDC_A(\mu, \nu) = \sup_{u \in A} \int u d(\mu - \nu):$$

where $A = \{u : u \text{ convex and } u \geq 0 \text{ almost everywhere}\}$.

Remark that since $0 \in K$, $VDC_A(\mu, \nu) \geq 0$ for all μ, ν because the zero function belongs to the set A .

Proposition

$$VDC_A(\mu, \nu) = 0 \iff \int u d\mu \geq \int u d\nu \text{ for all } u \in A$$

Intuition: $VDC_A(\mu, \nu) = 0 \iff \int u d\mu \geq \int u d\nu$ for all $u \in A$
 $\iff \int u d\mu \geq \int u d\nu$ for all u convex.

Variational Dominance Criterion (VDC)

Definition (Variational Dominance Criterion (VDC))

Given a bounded open convex subset $A \subset \mathbb{R}^d$, a pair of Borel probability measures $\mu, \nu \in \mathcal{P}(A)$, and a compact set $K \subset \mathbb{R}^d$ ($\emptyset \subset K$), define:

$$VDC_A(\mu, \nu; K) = \sup_{u \in \mathcal{U}_A} \int_A u d(\mu - \nu):$$

where $\mathcal{U}_A = \{u : \mathbb{R} \rightarrow \mathbb{R}; u \text{ convex and } u \geq 0 \text{ almost everywhere}\}$.

Remark that since $0 \in K$, $VDC_A(\mu, \nu; K) \geq 0$ for all μ, ν because the zero function belongs to the set \mathcal{U}_A .

Proposition

$$VDC_A(\mu, \nu; K) = 0 \iff \int_A u(x) d\mu(x) \leq \int_A u(x) d\nu(x) \text{ for all } u \in \mathcal{U}_A$$

Intuition: $VDC_A(\mu, \nu; K) = 0 \iff \int_A u(x) d\mu(x) \leq \int_A u(x) d\nu(x)$ for all $u \in \mathcal{U}_A$
 $\iff \int_A u(x) d\mu(x) \leq \int_A u(x) d\nu(x)$ for all u convex.

Informally, the proposition implies that $VDC_A(\mu, \nu; K)$ is small when μ is more spread out than ν , and large otherwise

Input Convex Maxout Networks

Problem: Statistical rates of estimation of the VDC are cursed by dimension, i.e. $\mathbb{E}[\text{VDC}_K(\hat{w}; j)] - \text{VDC}_K(w; j) \leq C n^{-2/d}$. The set of convex functions is too large (its Rademacher complexity scales like $n^{-2/d}$).

Input Convex Maxout Networks

Problem: Statistical rates of estimation of the VDC are cursed by dimension, i.e. $\mathbb{E}[\text{VDC}_K(\hat{\theta}; j)] - \text{VDC}_K(\theta; j) \leq C n^{-2/d}$. The set of convex functions is too large (its Rademacher complexity scales like $n^{-2/d}$).

Idea: **Approximate convex functions with bounded gradients using neural networks**

Input Convex Maxout Networks

Problem: Statistical rates of estimation of the VDC are cursed by dimension, i.e. $\mathbb{E}[\text{VDC}_K(\hat{\theta}; j)] \sim \text{VDC}_K(\theta; n; j) \cdot C n^{2=d}$. The set of convex functions is too large (its Rademacher complexity scales like $n^{2=d}$).

Idea: **Approximate convex functions with bounded gradients using neural networks**

Previous work: Input Convex Neural Networks (Amos et al., 2017). But we can do better in our setting!

Input Convex Maxout Networks

Problem: Statistical rates of estimation of the VDC are cursed by dimension, i.e. $\sqrt{\text{VDC}_K(\cdot)} \propto \sqrt{\text{VDC}_K(\cdot; n)} \propto \sqrt{n} \propto 2^{d/2}$. The set of convex functions is too large (its Rademacher complexity scales like $n^{-1/2}$).

Idea: **Approximate convex functions with bounded gradients using neural networks**

Previous work: Input Convex Neural Networks (Amos et al., 2017). But we can do better in our setting!

Idea: maximum of a few functions are good approximations of convex functions.

Can we stack them in layers? Yes! **Input Convex Maxout Networks.**

Figure 2: Shallow maxout network. ICMNs are maxout networks with convex increasing activations such that all weights beyond the first layer are non-negative.

$F_{L;M;k;+}(1)$: set of ICMNs with fixed architecture and bound on weights, such that $F_{L;M;k}(1) \subseteq A$.

$F_{L;M;k;+}(1)$: set of ICMNs with fixed architecture and bound on weights, such that $F_{L;M;k}(1) \subseteq A$.

We replace $A = \{u : \exists R; u \text{ convex and } r \leq K \text{ a.e.g by } F_{L;M;k}(1)$, and obtain the surrogate VDC:

$$\text{VDC}_{F_{L;M;k;+}(1)}(x, j) = \sup_{u \in F_{L;M;k;+}(1)} \inf_{z} u(z): \quad (2)$$

$F_{L;M;k;+}(1)$: set of ICMNs with fixed architecture and bound on weights, such that $F_{L;M;k}(1) \subseteq A$.

We replace $A = \{u \in \mathbb{R}^n; u \text{ convex and } \|u\|_2 \leq K \text{ a.e.g}\}$ by $F_{L;M;k}(1)$, and obtain the surrogate VDC:

$$\text{VDC}_{F_{L;M;k;+}(1)}(\alpha, \beta) = \sup_{u \in F_{L;M;k;+}(1)} \alpha u^T \beta \quad (2)$$

The surrogate VDC solves two problems at once:

It enjoys parametric estimation rates:

$$|\text{VDC}_{F_{L;M;k;+}(1)}(\alpha, \beta) - \text{VDC}_{F_{L;M;k;+}(1)}(\alpha, \beta; n)| \leq C n^{-1/2}.$$

We can use gradient descent to solve the variational problem (2) (no guarantees, but it works in practice).

We take a base generator g_0 trained using the baseline GAN training loss, and consider the problem:

$$\min_{g \in \mathcal{G}} \max_{f \in \mathcal{F}} E_{X \sim p_n} [f(X)] - E_{Y \sim p_0} [f(g(Y))] + \text{VDC}_{F; M; k; + (1)}(g_{\#} \parallel (g_0)_{\#}) \quad (3)$$

Here $g_{\#}$ is the distribution of the generated samples $g(X)$, $X \sim p_n$.

That is, we add the surrogate VDC between the learned and the baseline distribution: we want to bias $g_{\#}$ to be more spread-out than $(g_0)_{\#}$.

Mode collapse mitigation: mixture of Gaussians

The target r is a mixture of 8 gaussians in two dimensions

g_0 is a mode collapsed generator

g is trained with WGAN-GP penalized with the surrogate VDC.

Mode collapse mitigation: mixture of Gaussians

The target p_r is a mixture of 8 gaussians in two dimensions

g_0 is a mode collapsed generator

g is trained with WGAN-GP penalized with the surrogate VDC.

Figure 3: Probing mode collapse for GAN training. A converged generator needs to have a low negative likelihood and low mode collapse score. Collapse score: KL divergence between the discrete distribution obtained by assignment to closest centroid, and uniform distribution.

GAN experiments in high dimensions

Table 1: FID scores for WGAN-GP and WGAN-GP with VDC surrogate for convex functions approximated by either ICNNs with softplus activations or ICMNs, on the CIFAR-10 dataset. ICMNs improve upon the baseline g_0 and outperform ICNNs with softplus. FID score for WGAN-GP + VDC includes mean values \pm one standard deviation for 5 repeated runs with different random initialization seeds.

	FID	
g_0 : WGAN-GP	69.67	
g : WGAN-GP + VDC CP-Flow ICNN	83.470	3.732
g : WGAN-GP + VDC ICMN (Ours)	67.317	0.776

Portfolio optimization (Post et al., 2018; Xue et al., 2020): The goal is to find a portfolio G_2 that enhances a benchmark portfolio G_1 in a certain way: the return of G_2 must have high expectation, but its distribution must be less spread out than for G_1 — less risk

Distributional reinforcement learning (Martin et al., 2020): We want to learn policies with dominance constraints on the distribution of the reward.

Thank you!

Contacts: cd2754@nyu.edu, yzs2@cornell.edu , mroueh@us.ibm.com

- Amos, B., Xu, L., and Kolter, J. Z. (2017). Input convex neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 146–155. PMLR.
- Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein gan. *arXiv preprint arXiv:1701.07875*.
- Ekeland, I. and Schachermayer, W. (2014). Optimal transport and the geometry of $L^1(\mathbb{R}^d)$. *Proceedings of the American Mathematical Society*, 142.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680.
- Martin, J. D., Lyskawinski, M., Li, X., and Englot, B. (2020). Stochastically dominant distributional reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*. JMLR.org.
- Post, T., Karabati, S., and Arvanitis, S. (2018). Portfolio optimization based on stochastic dominance and empirical likelihood. *Journal of Econometrics*, 206(1):167–186.
- Xue, M., Shi, Y., and Sun, H. (2020). Portfolio optimization with relaxation of stochastic second order dominance constraints via conditional value at risk. *Journal of Industrial and Management Optimization*, 16(6):2581–2602.