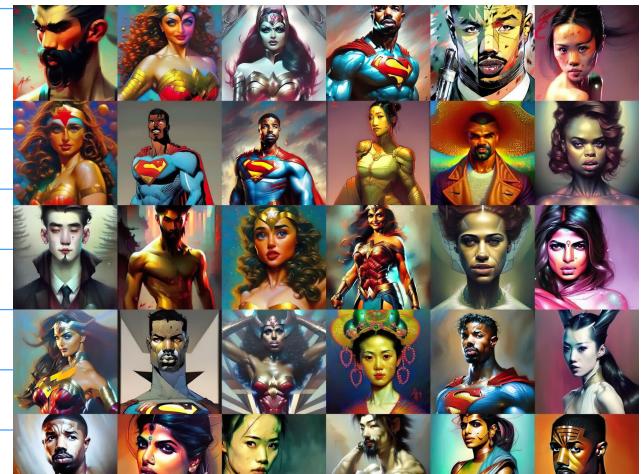


# Lecture 21 · Generative modeling with Kernels

Our goal in this lecture is to put together our current knowledge of kernel methods towards one of the modern problems of interest in Data Sci. call **Generative modeling**:

- AI art generators such as DALL-E & Stable diffusion
- chatbots
- AI voice generators
- etc.



Generative modeling is broadly an unsupervised learning task.

Goal: Given a set of samples  $\underline{y}_1, \dots, \underline{y}_n \stackrel{iid}{\sim} \mathcal{D}$  (target measure  $\mathcal{D}$ ) generate a new sample  $\underline{y}_{n+1} \sim \mathcal{D}$ .

why is it difficult? generally  $\mathcal{D} \in \mathcal{P}(X)$  when  $X$  is very high dim. (eg. high resolution images) & we only have access to a small set of samples of  $\mathcal{D}$ , i.e.,  $\underline{y}_1, \dots, \underline{y}_n$ .

Sampling is an old prob. but generative modeling is relatively new:

- Sampling is broadly the problem of generating iid samples from some measure  $\pi$ . Typically, we assume we knew  $\pi$  or at least its density up to a constant.

↳ Markov chain Monte Carlo

(Casella & Roberts, "Monte Carlo Statistical Methods")

↳ Variational inference

(Blei et al., "Variational Inference: A review for statisticians")

- Generative modeling is a sampling problem

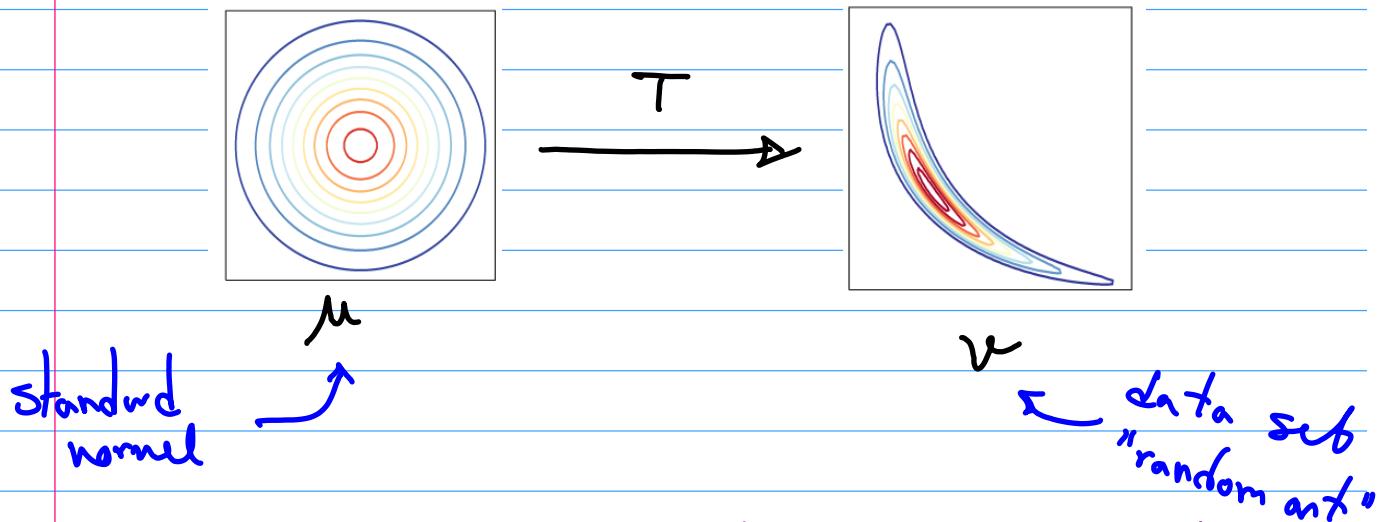
The novelty is that we only have samples from  $\pi$ !

In this lecture we will cast generative modeling as a transport problem & propose a solution using kernel methods.

Def<sup>n</sup>(Push forward) Given a prob. measure  $\mu \in \mathcal{P}(\mathcal{X})$  & a map  $T: \mathcal{X} \rightarrow \mathcal{Y}$  we define the push forward of  $\mu$  with  $T$ , denoted as  $T_{\#}\mu \in \mathcal{P}(\mathcal{Y})$ ) as  $(T_{\#}\mu)(A) := \mu(T^{-1}(A))$ ,  $\forall$  Borel sets  $A$ . Here  $T^{-1}(A) := \{x \in \mathcal{X} \text{ s.t. } T(x) \in A\}$  (the image of  $A$ )

Simple intuition for random variables:

If  $x \sim \mu$  then  $T(x) \sim T_{\#}\mu$ .



Transport is also an old problem & a "hot" math subject these days. It is the main topic in the field of Optimal Transport (Villani, "Optimal Transport, old & new").

## 21.1 Relaxation Using MMD

Problem: Given  $\mathcal{Y} = \{\underline{y}_1, \dots, \underline{y}_n\} \subset \mathbb{R}^d$  where  $\underline{y}_j \stackrel{iid}{\sim} \nu \in \mathcal{P}(\mathbb{R}^d)$  find a transport map  $\bar{T}: \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that

$$T_{\#}\mu \approx \nu \quad \text{where } \mu = N(0, I).$$

(3)

It is natural to formulate the above prob.  
as an optimization problem for  $T$ ! that is,

Find  $\bar{T}$  such that  $\bar{T} \# \mu$  is "close" to

$$v^n = \frac{1}{n} \sum_{j=1}^n \delta_{\underline{y}_j}, \text{ what do we mean by close?}$$

(II)

$$T = \underset{S: \mathbb{R}^d \rightarrow \mathbb{R}^d}{\operatorname{argmin}} \text{MMD}_K(S \# \mu, v^n)$$

where  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is an appropriate kernel such  
as RBF or RQ ( $\epsilon$  lengthscale should be fixed!).

Dealing with  $\bar{T} \# \mu$  is a bit cumbersome so let  
us replace  $\mu$  also with an empirical measure.

$$\mu^n = \frac{1}{n} \sum_{j=1}^n \delta_{\underline{x}_j}, \quad \underline{x}_j \sim M = N(0, I)$$

$$\Rightarrow T \# \mu \approx T \# \mu^n = \frac{1}{n} \sum_{j=1}^n \delta_{T(\underline{x}_j)}$$

Notation:  $\mathcal{Y} = \{\underline{y}_1, \dots, \underline{y}_n\}$

$$\mathcal{X} = \{\underline{x}_1, \dots, \underline{x}_n\}$$

$$T(\mathcal{X}) = \{T(\underline{x}_1), \dots, T(\underline{x}_n)\}$$

(4)

Farther approx. (T1) with

$$(T_2) \quad T = \underset{S: \mathbb{R}^d \rightarrow \mathbb{R}^d}{\operatorname{argmin}} \text{MMD}_K(S(X), Y)$$

Observe, this is now a nice loss since we can easily evaluate it for any choice of S!

But, the map T can be anything! in fact, prob. has infinitely many solutions!

## 21.2 Regularization with RKHS norms & a Rep.

### Formula

We would like to bring some stability to the problem by regularizing T. Also, in practice, having a T that is "nicer/more regular" would be beneficial. So, we may decide to penalize an RKHS norm of T.

Challenge:  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$  how do we make sense of its RKHS norm?

Simple solution: Simply add up the RKHS norms of components of the map

$$T(\underline{x}) = (\underline{T}^1(\underline{x}), \underline{T}^2(\underline{x}), \dots, \underline{T}^d(\underline{x}))^T \in \mathbb{R}^d$$

(3)

$$\|T\|_Q^2 := \sum_{j=1}^d \|T^j\|_Q^2$$

for some kernel  $Q: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , e.g. RBF, Matérn, Polynomial, etc.

We can now define our regularized transport problem:

$$(T_3) T = \underset{S: \mathbb{R}^d \rightarrow \mathbb{R}^d}{\operatorname{argmin}} \text{MMD}_K(S(X), Y) + \lambda \sum_{j=1}^d \|S^j\|_Q^2$$

**WARNING:** The kernels  $K$  &  $Q$  are not necessarily related! They need to be chosen & tuned separately.

Observe, Prob. (T3) can also be viewed as a "regression" prob. for the components  $T^j$  of  $T$ !  
So we can write a rep. formula for it!

Let us write  $S(\underline{x}_i) = \underline{z}_i$ ,  $i=1, \dots, n$

&  $\underline{z} = S(X)$  then we can write (T3)  
as the new prob.

(6)

$$\min_{\underline{z} = \{\underline{z}_1, \dots, \underline{z}_n\}} \min_{S: \mathbb{R}^d \rightarrow \mathbb{R}^d} MMD_K(\underline{z}, \gamma) + \lambda \sum_{j=1}^d \|S^j\|_Q^2$$

equivalently  $\rightarrow$

$\text{s.t. } S(\underline{x}_i) = \underline{z}_i, i=1, \dots, n$

$$S^j(\underline{x}_i) = z_{ij}, i=1, \dots, n$$

$$j=1, \dots, d$$

Minimizing over  $S$  for a fixed  $\underline{z}$  we get that

$$S^j = Q(\cdot, X) Q(X, X)^{-1} (\underline{z}_{:,j})^T$$

& further more  $\|S^j\|_Q^2 = z_{:,j} Q(X, X)^{-1} (\underline{z}_{:,j})^T$   
 so the rep. then takes the form

$$T(\underline{y}) = (T^1(\underline{x}), \dots, T^d(\underline{x}))^T$$

$$T^j = Q(\cdot, X) Q(X, X)^{-1} (\underline{z}_{:,j})^T, j=1, \dots, d$$

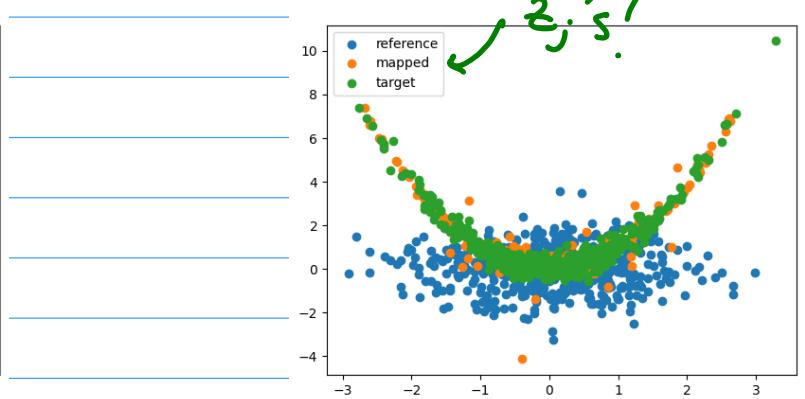
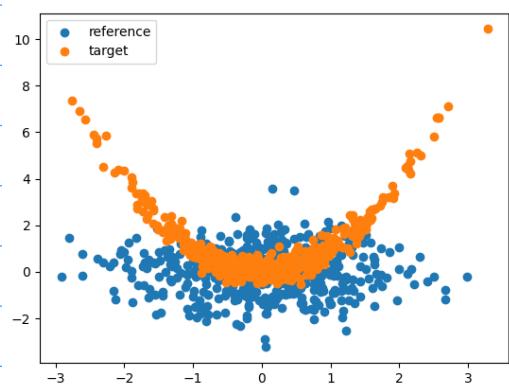
where  $\underline{z} \in \mathbb{R}^{n \times d}$  solves

$$\min_{\underline{z}} MMD_K(\underline{z}, \gamma) + \lambda \sum_{j=1}^d z_{:,j} Q(X, X)^{-1} (\underline{z}_{:,j})^T$$

The above Prob. is easily solvable using off-the shelf optimizers.

(7)

mapping a 2D Gaussian to a non gaussian measure



The code uses a RF implementation of MMD for speed. Also looks for maps of the form

$\underline{x} \mapsto \underline{x} + T(\underline{x})$  rather than  $\underline{x} \mapsto T(\underline{x})$



