

A Mean-Field Theory for Kernel Learning in Generative and Discriminative Models of Machine Learning

MASOUD BADIEI KHUZANI

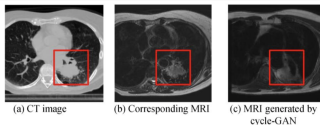
NOVEMBER, 2019

STANFORD UNIVERSITY

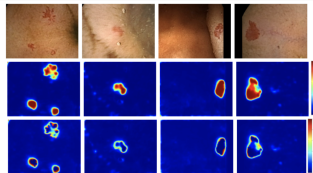
Machine Learning Applications in Medical Imaging

- Machine learning and AI techniques have changed the landscape of medical imaging.
- Many medical tasks such as segmentation, registration, or diagnosis can be done automatically without any intervention from clinicians.

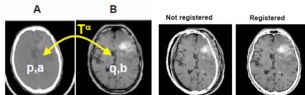
1. Modality Transformation



2. Segmentation



3. Registration



Machine Learning Problems in a Nutshell

-A broad range of machine learning tasks can be reduced to the problem of learning a target function

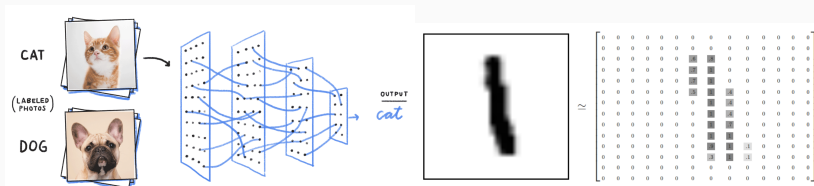
$$h : \mathcal{X} \rightarrow \mathcal{Y}.$$

Machine Learning Problems in a Nutshell

-A broad range of machine learning tasks can be reduced to the problem of learning a target function

$$h : \mathcal{X} \rightarrow \mathcal{Y}.$$

- **Discriminative (classification) model:** $\mathcal{X} = \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n}, \mathcal{Y} = \{0, 1\},$
 $\mathcal{S} = \{0, 1, \dots, 255\}.$

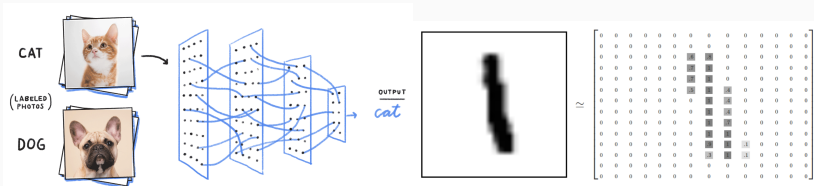


Machine Learning Problems in a Nutshell

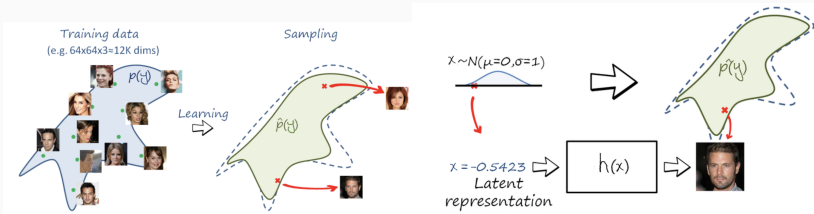
-A broad range of machine learning tasks can be reduced to the problem of learning a target function

$$h: \mathcal{X} \rightarrow \mathcal{Y}.$$

- **Discriminative (classification) model:** $\mathcal{X} = \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n}, \mathcal{Y} = \{0, 1\},$
 $\mathcal{S} = \{0, 1, \dots, 255\}.$



- **Generative (sampling) model:** $\mathcal{X} = \mathbb{R}, \mathcal{Y} = \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n} \times \mathcal{S}^{n \times n}$



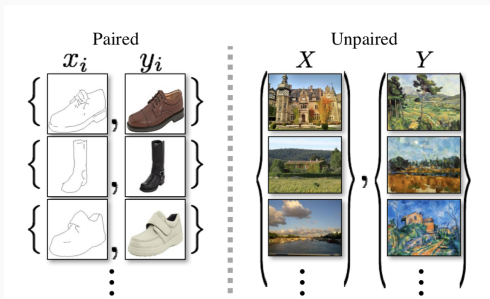
Supervised vs Unsupervised Models

- **Supervised models:** The paired instances from the target map are available

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, \quad y_i = h(x_i). \quad (1)$$

- **Unsupervised models:** The unpaired instances from each domain is available (for classification it means class label is unavailable)

$$x_1, \dots, x_n \in \mathcal{X}, \quad y_1, \dots, y_m \in \mathcal{Y}, \quad y_i \neq h(x_i).$$



- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} D \left(\hat{P}^n = \frac{1}{n} \sum_{i=1}^n \delta(y - f(x_i)) \parallel \hat{Q}^m = \frac{1}{m} \sum_{i=1}^m \delta(y - y_i) \right),$$

$D(\cdot, \cdot)$: divergence between two distributions, $\delta(\cdot)$: Dirac's delta function.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} D \left(\hat{P}^n = \frac{1}{n} \sum_{i=1}^n \delta(y - f(x_i)) \parallel \hat{Q}^m = \frac{1}{m} \sum_{i=1}^m \delta(y - y_i) \right),$$

$D(\cdot, \cdot)$: divergence between two distributions, $\delta(\cdot)$: Dirac's delta function.

- The variational form of the distance between two distributions P, Q defined on \mathcal{X} :

$$D(P \parallel Q) = \max_{g \in \mathcal{G}} \left| \int_{\mathcal{X}} g(x) dP(x) - \int_{\mathcal{X}} g(x) dQ(x) \right|. \quad (2)$$

\mathcal{G} : function class.

Variational form of the divergence

- **Total Variation Distance:**

$$\text{TV}(P||Q) = \sup_{g:\mathcal{X}\rightarrow\mathbb{R}:\|g\|_{\infty}\leq 1/2} \left| \int_{\mathcal{X}} g(x)dP(x) - \int_{\mathcal{X}} g(x)dQ(x) \right|.$$

- **Wasserstein (KantorovichRubinstein) Distance:**

$$W_1(P||Q) = \sup_{g:\mathcal{X}\rightarrow\mathbb{R}:\text{Lip}(g)\leq 1} \left| \int_{\mathcal{X}} g(x)dP(x) - \int_{\mathcal{X}} g(x)dQ(x) \right|.$$

- **Maximum Mean Discrepancy Distance:**

$$\text{MMD}_{\mathcal{H}}(P||Q) = \sup_{g:\mathcal{X}\rightarrow\mathbb{R}:\|g\|_{\mathcal{H}}\leq 1} \left| \int_{\mathcal{X}} g(x)dP(x) - \int_{\mathcal{X}} g(x)dQ(x) \right|.$$

How to Learn The Target Map $h : \mathcal{X} \rightarrow \mathcal{Y}$

*

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

Plug \rightarrow

$$\min_{f \in \mathcal{F}} D \left(\hat{P}^n = \frac{1}{n} \sum_{i=1}^n \delta(y - f(x_i)) \parallel \hat{Q}^m = \frac{1}{m} \sum_{i=1}^m \delta(y - y_i) \right),$$

$D(\cdot, \cdot)$: divergence between two distributions, $\delta(\cdot)$: Dirac's delta function.

- The variational form of the distance between two distributions P, Q defined on \mathcal{X} :

$$D(P \parallel Q) = \max_{g \in \mathcal{G}} \left| \int_{\mathcal{X}} g(x) dP(x) - \int_{\mathcal{X}} g(x) dQ(x) \right|. \quad (3)$$

\mathcal{G} : function class.

How to Learn The Target Map $h : \mathcal{X} \rightarrow \mathcal{Y}$?



- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m g(y_i) - \frac{1}{n} \sum_{i=1}^n g(f(x_i)) \right|.$$

\mathcal{G} : function class.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

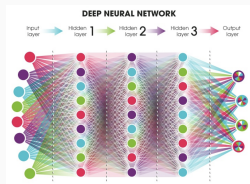
\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m g(y_i) - \frac{1}{n} \sum_{i=1}^n g(f(x_i)) \right|.$$

\mathcal{G} : function class.

Deep Neural Networks



-In deep neural networks with n -layers, the parametric function takes the following form

$$\text{NN}(x; \mathbf{W}) \stackrel{\text{def}}{=} \sigma(W_n \sigma(\dots \sigma(W_1 x))), \quad W_k \in \mathbb{R}^{m_k \times m_{k-1}},$$

where $\mathbf{W} \stackrel{\text{def}}{=} (W_1, \dots, W_n)$.

$-\sigma(\cdot)$ is the activation function for non-linear function approximation.

1. ReLU: $\sigma(x) = \max\{0, x\}$.
2. logisitc: $\sigma(x) = \frac{1}{1+e^{-x}}$.
3. arctan: $\sigma(x) = \tan^{-1}(x)$.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m g(y_i) - \frac{1}{n} \sum_{i=1}^n g(f(x_i)) \right|.$$

\mathcal{G} : function class.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_n} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \text{NN}(x_i, \mathbf{W})) + \sum_{i=1}^n \frac{\lambda_i}{2} \|\mathbf{W}_i\|_F^2.$$

$$\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}, x \mapsto f(x) = \text{NN}(x_i, \mathbf{W}), \mathbf{W} \in \prod_{k=1}^n \mathbb{R}^{m_k \times m_{k-1}}\}.$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{\mathbf{W}_1, \dots, \mathbf{W}_n} \max_{\mathbf{W}'_1, \dots, \mathbf{W}'_n} \left| \frac{1}{m} \sum_{i=1}^m \text{NN}(\mathbf{W}'_i, y_i) - \frac{1}{n} \sum_{i=1}^n \text{NN}(\mathbf{W}'_i, \text{NN}(\mathbf{W}, x_i)) \right|.$$

$$\mathcal{G} = \{g : \mathcal{Y} \rightarrow \mathbb{R}, x \mapsto g(y) = \text{NN}(y, \mathbf{W}), \mathbf{W} \in \prod_{k=1}^n \mathbb{R}^{m_k \times m_{k-1}}\}.$$

How to Learn The Target Map $h : \mathcal{X} \rightarrow \mathcal{Y}$?

*

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_n} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \text{NN}(x_i, \mathbf{w})) + \sum_{i=1}^n \frac{\lambda_i}{2} \|\mathbf{w}_i\|_F^2.$$

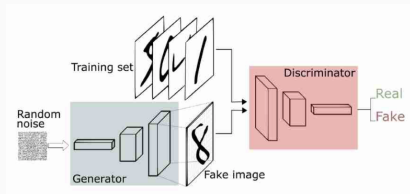
$$\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y}, x \mapsto f(x) = \text{NN}(x_i, \mathbf{w}), \mathbf{w} \in \prod_{k=1}^n \mathbb{R}^{m_k \times m_{k-1}}\}.$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_n} \max_{\mathbf{w}'_1, \dots, \mathbf{w}'_n} \left| \frac{1}{m} \sum_{i=1}^m \text{NN}(\mathbf{w}', y_i) - \frac{1}{n} \sum_{i=1}^n \text{NN}(\mathbf{w}', \text{NN}(\mathbf{w}, x)) \right|.$$

Generative Adversarial Networks

$$\mathcal{G} = \{g : \mathcal{Y} \rightarrow \mathbb{R}, x \mapsto g(y) = \text{NN}(y, \mathbf{w}), \mathbf{w} \in \prod_{k=1}^n \mathbb{R}^{m_k \times m_{k-1}}\}.$$



Kernel Machines

- Let \mathcal{H} denotes the Hilbert space (inner product space with Cauchy sequence limits) of real-valued functions on \mathcal{X} .
- For $x \in \mathcal{X}$, consider the map $L_x : \mathcal{H} \rightarrow \mathbb{R}, f \mapsto L_x[f] = f(x)$. If L_x is a bounded operator, we say \mathcal{H} is a reproducing kernel Hilbert space (RKHS).
- $L_x \in \mathcal{H}^*$, where \mathcal{H}^* is the dual-space of the Hilbert space \mathcal{H} .
- By Riesz representation theorem, there exists an element $\phi(x) \in \mathcal{H}$, such that

$$L_x(f) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \quad \forall f \in \mathcal{H}.$$

- In particular, $\phi : \mathcal{X} \rightarrow \mathcal{H}$ and $\phi(y) \in \mathcal{H}$. Therefore,

$$L_x(\phi(y)) = \langle \phi(y), \phi(x) \rangle_{\mathcal{H}} \stackrel{\text{def}}{=} K(x, y).$$

- $K(x, y)$: kernel function, $\phi(x)$: feature map.



Moore-Aronszajn Theorem . Consider the kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is symmetric and is positive definite in the sense that

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(x_i, x_j) > 0,$$

for all $n \in \mathbb{N}$, $c_1, \dots, c_n \in \mathbb{R}$, and $x_1, \dots, x_n \in \mathcal{X}$. Then, there exists a unique Hilbert space \mathcal{H}_K for which K has the reproducing property. Furthermore, for every function $f \in \mathcal{H}_K$, we have the following expansion

$$f(x) = \sum_{i=1}^{\infty} w_i K(x, x_i), \quad \text{for some } x_1, x_2, \dots \in \mathcal{X}.$$

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), f(x_i)).$$

\mathcal{F} : a function class, ℓ : a loss function.

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{i=1}^m g(y_i) - \frac{1}{n} \sum_{i=1}^n g(f(x_i)) \right|.$$

\mathcal{G} : function class.

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set

$$\min_{(w_1, \dots, w_n) \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell\left(y_i, \sum_{j=1}^n w_j K(x_i, x_j)\right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

$$\mathcal{F} = \mathcal{H}_K = \left\{ f : \mathcal{X} \rightarrow \mathcal{Y}, x \mapsto \sum_{i=1}^n w_i K(x, x_i), \mathbf{w} \in \mathbb{R}^n \right\}$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{w_1, \dots, w_n} \max_{\|g\|_{\mathcal{H}_K} \leq 1} \left| \frac{1}{m} \sum_{i=1}^m g(y_i) - \frac{1}{n} \sum_{i=1}^n g(\text{NN}(\mathbf{W}, x_i)) \right|.$$

$$\mathcal{G} = \{g : \mathcal{X} \rightarrow \mathbb{R} : \|g\|_{\mathcal{H}_K} \leq 1\}.$$

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set ($\ell(y, x) = \max\{0, 1 - xy\}$.)

$$\min_{(w_1, \dots, w_n) \in \mathbb{R}^n} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell\left(y_i, \sum_{j=1}^n w_j K(x_i, x_j)\right)}_{\text{Kernel Support Vector Machines}} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{w_1, \dots, w_n} \text{MMD}_K(\hat{P}_W^n \| \hat{Q}^m) = \frac{1}{m^2} \sum_{i,j=1}^m K(y_i, y_j) + \frac{1}{n^2} \sum_{i,j=1}^n K(\text{NN}(\mathbf{W}, x_i), \text{NN}(\mathbf{W}, x_j)) - \underbrace{\frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m K(y_j, \text{NN}(\mathbf{W}, x_i))}_{\text{Generative Moment Matching Networks}}$$

, where $y_1, \dots, y_m \sim_{\text{i.i.d.}} Q$, and $\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n) \sim_{\text{i.i.d.}} P_W$.

Kernel Selection Problem

Kernel Model Selection Problem: A bad kernel may yield a poor machine learning system.

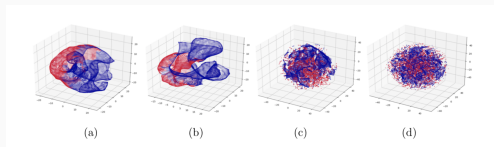


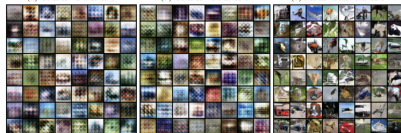
Figure 1: t -SNE plot for feature maps generated by Gaussian kernel $K(x, y) = e^{-\gamma \|x-y\|_2}$ for different bandwidth values $\gamma > 0$.



(a) GMM-D MNIST

(b) GMM-C MNIST

(c) MMD GAN MNIST



(d) GMM-D CIFAR-10

(e) GMM-C CIFAR-10

(f) MMD GAN CIFAR-10

Kernel Optimization

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set ($\ell(y, x) = \max\{0, 1 - xy\}$.)

$$\min_{(w_1, \dots, w_n) \in \mathbb{R}^n} \max_{K \in \mathcal{K}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell\left(y_i, \sum_{j=1}^n w_j K(x_i, x_j)\right)}_{\text{Kernel Support Vector Machines}} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\min_{w_1, \dots, w_n} \max_{K \in \mathcal{K}} \text{MMD}_K(\hat{P}_{\mathbf{W}}^n \|\hat{Q}^m) = \frac{1}{m^2} \sum_{i,j=1}^m K(y_i, y_j) + \frac{1}{n^2} \sum_{i,j=1}^n K(\text{NN}(\mathbf{W}, x_i), \text{NN}(\mathbf{W}, x_j)) \\ - \underbrace{\frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m K(y_j, \text{NN}(\mathbf{W}, x_i))}_{\text{Generative Moment Matching Networks}}$$

, where $y_1, \dots, y_m \sim_{\text{i.i.d.}} Q$, and $\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n) \sim_{\text{i.i.d.}} P_{\mathbf{W}}$.

- Suppose the kernel function is shift invariant $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{x} - \mathbf{y})$.
- Rahimi and Recht ¹ proved that shift invariant kernels have the following Fourier representation

$$K(\mathbf{x} - \mathbf{y}) = \mathbb{E}_{\mu}[\varphi(\mathbf{x}, \boldsymbol{\xi})\varphi(\mathbf{y}, \boldsymbol{\xi})] = \int_{\mathbb{R}^D} \varphi(\mathbf{x}, \boldsymbol{\xi})\varphi(\mathbf{y}, \boldsymbol{\xi})d\mu(\boldsymbol{\xi}),$$

where $\varphi(\mathbf{x}, \boldsymbol{\xi}) = \cos(\langle \mathbf{x}, \boldsymbol{\xi} \rangle + b)$, $b \sim \text{Uniform}([0, 2\pi])$.

- $\varphi(\mathbf{x}, \boldsymbol{\xi})$ is called the random feature map since

$$K(\mathbf{x} - \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{H}} = \langle \varphi(\mathbf{x}, \boldsymbol{\xi}), \varphi(\mathbf{y}, \boldsymbol{\xi}) \rangle_{L^2(\mu)}.$$

¹Ali Rahimi and Benjamin Recht. *Random Features for Large-Scale Kernel Machines*
NIPS 2007.

Kernel Optimization

- **Supervised models:** the target map is learned by minimizing a risk function over a given training set ($\ell(y, x) = \max\{0, 1 - xy\}$.)

$$\min_{(w_1, \dots, w_n) \in \mathbb{R}^n} \max_{\mu \in \mathcal{P}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell\left(y_i, \sum_{j=1}^n w_j \mathbb{E}_{\mu}[\varphi(x_i, \boldsymbol{\xi}) \varphi(x_j, \boldsymbol{\xi})]\right)}_{\text{Kernel Support Vector Machines}} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- **Unsupervised models:** the target map is learned by minimizing the distance between the empirical distribution of data and the model output

$$\begin{aligned} \min_{W_1, \dots, W_n} \max_{\mu \in \mathcal{P}} \text{MMD}_K(\hat{P}_W^n \| \hat{Q}^m) &= \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}_{\mu}[\varphi(y_i, \boldsymbol{\xi}) \varphi(y_j, \boldsymbol{\xi})] \\ &+ \frac{1}{n^2} \sum_{i,j=1}^n \mathbb{E}_{\mu}[\varphi(\text{NN}(\mathbf{W}, x_i), \boldsymbol{\xi}) \varphi(\text{NN}(\mathbf{W}, x_i), \boldsymbol{\xi})] \\ &- \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m \mathbb{E}_{\mu}[\varphi(y_j, \boldsymbol{\xi}) \varphi(\text{NN}(\mathbf{W}, x_i), \boldsymbol{\xi})]. \end{aligned}$$

, where $y_1, \dots, y_m \sim \text{i.i.d. } Q$, and $\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n) \sim \text{i.i.d. } P_W$.

- **Supervised Models:** For binary classification $\mathcal{Y} = \{-1, +1\}$ using kernel SVMs, we optimize the kernel target alignment.

$$\max_{\mu \in \mathcal{P}} \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} y_i y_j \mathbb{E}_{\mu}[\varphi(x_i, \xi) \varphi(x_j, \xi)].$$

- **Unsupervised Models:** We optimize the kernel target alignment ($M = m + n$)

$$\min_{W_1, \dots, W_n} \max_{\mu \in \mathcal{P}} \frac{2}{M(M-1)} \sum_{1 \leq i < j \leq n+m} z_i z_j \mathbb{E}_{\mu}[\varphi(v_i, \xi) \varphi(v_j, \xi)],$$

where $z_1, \dots, z_M \sim \text{Uniform}\{-1, +1\}$ and

1. If $z_i = +1$, let $v_i \in \{y_1, \dots, y_m\}$ (true data).
2. If $z_i = -1$ let $v_i \in \{\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n)\}$ (generated data).

- **Step 1:** Lets focus on the following optimization problem

$$\min_{W_1, \dots, W_n} \max_{\mu \in \mathcal{P}} \frac{2}{M(M-1)} \sum_{1 \leq i < j \leq n+m} z_i z_j \mathbb{E}_{\mu} [\varphi(v_i, \xi) \varphi(v_j, \xi)],$$

where $z_1, \dots, z_M \sim \text{Uniform}\{-1, +1\}$ and

1. If $z_i = +1$, let $v_i \in \{y_1, \dots, y_m\}$.
 2. If $z_i = -1$ let $v_i \in \{\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n)\}$.
- We rewrite the problem as a risk minimization

$$\min_{W_1, \dots, W_n} \max_{\mu \in \mathcal{P}} \frac{2}{\alpha M(M-1)} \sum_{1 \leq i < j \leq n} (\alpha z_i z_j - \mathbb{E}_{\mu} [\varphi(v_i, \xi) \varphi(v_j, \xi)])^2.$$

Kernel Optimization

- **Step 2:** We apply the Monte-Carlo sample average approximation.
- In particular, we optimize with respect to samples of the i.i.d. samples of the target distribution $\xi^1, \dots, \xi^N \sim_{\text{i.i.d.}} \mu$,

$$\min_{w_1, \dots, w_n} \max_{\hat{\mu}^N \in \mathcal{P}^N} \frac{2}{\alpha M(M-1)} \sum_{1 \leq i < j \leq n} \left(\alpha z_i z_j - \frac{1}{N} \sum_{k=1}^N \varphi(x, \xi_k) \varphi(y, \xi_k) \right)^2,$$

where $\mathcal{P}^N \stackrel{\text{def}}{=} \{\hat{\mu}^N \in \mathcal{M}(\mathbb{R}^D) : D(\hat{\mu}^N \| \hat{\mu}_0^N) \leq R\}$.

- We define the empirical distribution of the samples (particles) as below

$$\hat{\mu}^N(\xi) = \frac{1}{N} \sum_{i=1}^N \delta(\xi - \xi^i). \quad (4)$$

Particle Stochastic Gradient Descent

- **Step 3:** Use a particle stochastic gradient descent method to solve the risk minimization as follows:

1. Initialize the samples $\xi_0^1, \dots, \xi_0^N \sim \text{i.i.d. } \mu_0$.
2. At iteration $\ell = 0, 1, 2, \dots$, we draw two fresh labels $z_\ell, \tilde{z}_\ell \sim \text{Uniform}\{-1, +1\}$.
3. Sample $v_\ell \in \{y_1, \dots, y_m\}$ if $z_\ell = 1$, and $v_\ell \in \{\text{NN}(\mathbf{W}, x_1), \dots, \text{NN}(\mathbf{W}, x_n)\}$ if $z_\ell = -1$. We pick \tilde{v}_ℓ using a similar rule.
4. Apply the particle SGD with the step size $\eta > 0$

$$\xi_{\ell+1}^k = \xi_\ell^k - \frac{\eta}{N} \left(z_\ell \tilde{z}_\ell - \frac{1}{\alpha N} \sum_{k=1}^N \varphi(v_\ell; \xi_\ell^k) \varphi(\tilde{v}_\ell; \xi_\ell^k) \right) \nabla_{\xi} \left(\varphi(v_\ell; \xi_\ell^k) \varphi(\tilde{v}_\ell; \xi_\ell^k) \right),$$

for $k = 1, 2, \dots, N$.

5. Approximate the kernel

$$K_{\ell+1}(x, y) \approx \frac{1}{N} \sum_{k=1}^N \varphi(x; \xi_{\ell+1}^k) \varphi(y; \xi_{\ell+1}^k). \quad (5)$$

Evolution of the Histogram of SGD Particles

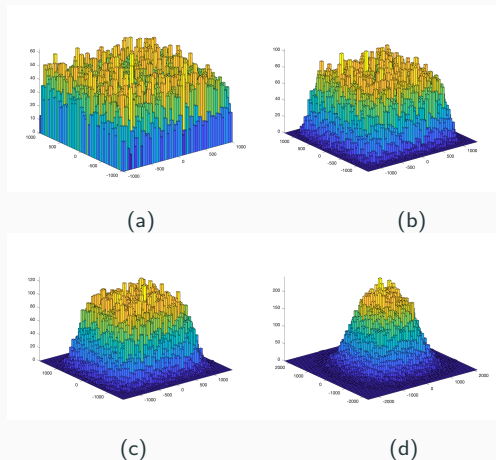


Figure 2: The evolution of the empirical measure $\mu_\ell^N(\boldsymbol{\xi}) = \frac{1}{N} \sum_{k=1}^N \delta(\boldsymbol{\xi} - \boldsymbol{\xi}_\ell^k)$ of the SGD particles $\boldsymbol{\xi}_\ell^1, \dots, \boldsymbol{\xi}_\ell^N \in \mathbb{R}^2$ at different iterations ℓ . The empirical measure of random feature maps seemingly converges to a Gaussian stationary measure corresponding to a Gaussian RBF kernel. Panel (a): $\ell = 0$, Panel (b): $\ell = 300$, Panel (c): $\ell = 1000$, and Panel (d): $\ell = 2500$.

Consistency of Monte-Carlo Approximations

Theorem: Consider the distribution ball with respect to the 2-Wasserstein distance $\mathcal{P} = \{\mu \in \mathcal{M}(\mathbb{R}^D) : W_2(\mu || \mu_0) \leq R\}$, where μ_0 is a user-defined distribution. Furthermore, consider

$$(\mathbf{W}_*, \mu_*) \stackrel{\text{def}}{=} \arg \min_{\mathbf{W} \in \mathcal{W}} \arg \sup_{\mu \in \mathcal{P}} \text{MMD}_\mu[P_{\mathbf{W}}, Q]$$
$$(\widehat{\mathbf{W}}_*^N, \widehat{\mu}_*^N) \stackrel{\text{def}}{=} \arg \min_{\mathbf{W} \in \mathcal{W}} \arg \inf_{\widehat{\mu}^N \in \mathcal{P}_N} \widehat{\text{MMD}}_{\widehat{\mu}^N}^\alpha[\widehat{P}_{\mathbf{W}}^n, \widehat{Q}^n],$$

where

$$\text{MMD}_\mu[P_{\mathbf{W}}, Q] = \mathbb{E}_{P \otimes 2}[\mathbb{E}_\mu[\varphi(\text{NN}(\mathbf{x}, \mathbf{W}), \boldsymbol{\xi})\varphi(\text{NN}(\mathbf{x}, \mathbf{W}), \boldsymbol{\xi})]] + \mathbb{E}_{Q \otimes 2}[\mathbb{E}_\mu[\varphi(y, \boldsymbol{\xi})\varphi(y', \boldsymbol{\xi})]] - 2\mathbb{E}_{P, Q}[\mathbb{E}_\mu[\varphi(\text{NN}(\mathbf{x}, \mathbf{W}), \boldsymbol{\xi})\varphi(y, \boldsymbol{\xi})]],$$

and

$$\widehat{\text{MMD}}_{\widehat{\mu}^N}^\alpha[\widehat{P}_{\mathbf{W}}^n, \widehat{Q}^n]$$
$$= \frac{2}{\alpha 2n(2n-1)} \sum_{1 \leq i < j \leq n} \left(\alpha z_i z_j - \frac{1}{N} \sum_{k=1}^N \varphi(\text{NN}(v_i, \mathbf{W}), \boldsymbol{\xi}_k) \varphi(\text{NN}(v_j, \mathbf{W}), \boldsymbol{\xi}_k) \right)^2.$$

Consistency of Monte-Carlo Approximations

Then, with the probability of (at least) $1 - 3\varrho$ over the training data samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and the random feature samples $\{\xi_0^k\}_{k=1}^N$, the following non-asymptotic bound holds

$$\begin{aligned} \left| \text{MMD}_\mu[P_{\mathbf{W}}, Q] - \widehat{\text{MMD}}_{\hat{\mu}^N}^\alpha[\hat{P}_{\mathbf{W}}^n, \hat{Q}^n] \right| \leq & \sqrt{\frac{L^2(d+2)}{N} \ln^{\frac{1}{2}} \left(\frac{2^8 N \text{diam}^2(\mathcal{X})}{\varrho} \right)} + \frac{8L^2}{\alpha} \\ & + 2 \max \left\{ \frac{c_1 L^2}{n} \ln^{\frac{1}{2}} \left(\frac{4}{\varrho} \right), \frac{c_2 R L^4}{n^2} \ln \left(\frac{4e \frac{L^4}{9}}{\varrho} \right) \right\}, \end{aligned}$$

where $c_1 = 3^{\frac{1}{4}} \times 2^4$, and $c_2 = 9 \times 2^{11}$.

Consistency of Monte-Carlo Approximations

Then, with the probability of (at least) $1 - 3\rho$ over the training data samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and the random feature samples $\{\xi_0^k\}_{k=1}^N$, the following non-asymptotic bound holds

$$\left| \text{MMD}_\mu[P_{\mathbf{W}}, Q] - \widehat{\text{MMD}}_{\hat{\mu}^N}^\alpha[\hat{P}_{\mathbf{W}}^n, \hat{Q}^n] \right| \leq \sqrt{\frac{L^2(d+2)}{N} \ln^{\frac{1}{2}} \left(\frac{2^8 N \text{diam}^2(\mathcal{X})}{\rho} \right)} + \frac{8L^2}{\alpha} + 2 \max \left\{ \frac{c_1 L^2}{n} \ln^{\frac{1}{2}} \left(\frac{4}{\rho} \right), \frac{c_2 R L^4}{n^2} \ln \left(\frac{4e \frac{L^4}{9}}{\rho} \right) \right\},$$

where $c_1 = 3^{\frac{1}{4}} \times 2^4$, and $c_2 = 9 \times 2^{11}$.

N : The number of random feature samples (Particles) ξ^1, \dots, ξ^N in Particle SGD.

n : The number of training samples $y_1, \dots, y_n \in \mathcal{Y}$, and $x_1, \dots, x_n \in \mathcal{X}$.

α : Regularization parameter of the empirical risk minimization.

Convergence of the empirical measure to a limiting measure

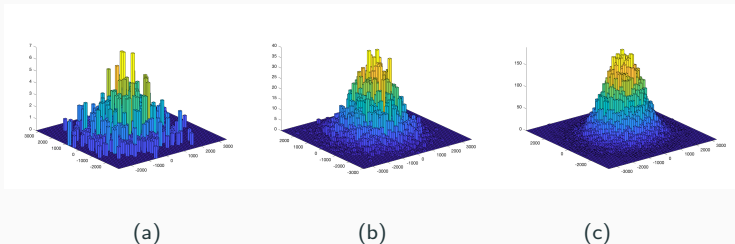


Figure 3: The histogram of the SGD particles at a fixed iteration $\ell = 10000$ and for different number of particles. Panel (a): $N = 1000$, Panel (b): $N = 10000$, Panel (c): $N = 50000$.

- As $N \rightarrow \infty$, it seems that $\hat{\mu}_\ell^N = \frac{1}{N} \sum_{k=1}^N \delta(\xi - \xi_\ell^k)$ converges to a limiting measure μ_ℓ^* .

Consistency of the Particle SGD

Theorem: (MCKEAN-VLASOV MEAN-FIELD PDE) Define the scaled empirical measure of the SGD particles embedded in the continuous time

$$\mu_t^N = \widehat{\mu}_{[Nt]}^N = \frac{1}{N} \sum_{k=1}^N \delta(\xi - \xi_{[Nt]}), \quad 0 \leq t \leq T.$$

Further, suppose that the Lebesgue density of the initial measure of particles $q_0(\xi) = \mu_0(d\xi)/d\xi$ exists. Then, there exists a unique solution $(p_t^*(\xi))_{0 \leq t \leq T}$ to the following non-linear partial differential equation

$$\begin{aligned} & \frac{\partial p_t(\xi)}{\partial t} \\ &= -\frac{\eta}{\alpha} \iint_{\mathcal{X} \times \mathcal{Y}} \left(\int_{\mathbb{R}^p} \varphi(v, \tilde{\xi}) \varphi(\tilde{v}, \tilde{\xi}) p_t(\tilde{\xi}) d\tilde{\xi} - \alpha z \tilde{z} \right) \nabla_{\xi}(p_t(\xi)) \nabla_{\xi}(\varphi(v; \xi) \varphi(\tilde{v}; \xi)) dP_{V,Z}^{\otimes 2}, \end{aligned} \tag{7}$$

where $P_{V,Z}$ has the marginals $P_{V|Z=+1} = P_W$ and $P_{V|Z=-1} = Q$. The PDE is initialized at $p_0(\xi) = q_0(\xi)$. Moreover, the measure-valued process $\{(\mu_t^N)_{0 \leq t \leq T}\}_{N \in \mathbb{N}}$ converges (weakly) to the unique solution $\mu_t^*(\xi) = p_t^*(\xi) d\xi$ as the number of particles tend to infinity $N \rightarrow \infty$.

Consistency of the Particle SGD

- The Mean-Field PDE can be viewed as the gradient flow for minimizing an energy functional

$$\frac{dp_t(\xi)}{dt} = -\eta \cdot \text{grad}_{p_t} E_\alpha(p_t(\xi)), \quad p_0(\xi) = q_0(\xi),$$

where $\text{grad}_{p_t} E(p_t(\xi)) = \nabla_\xi \cdot (p_t(\xi) \nabla_\xi R_\beta(p_t(\xi)))$ is the Riemannian gradient of $R_\beta(\mu_t(\xi))$ with respect to the metric of the Wasserstein manifold, and

$$\begin{aligned} \inf_{\mu \in \mathcal{M}(\mathbb{R}^p)} E_\alpha(p_t(\xi)) &\stackrel{\text{def}}{=} \frac{1}{\alpha} \int_{\mathbb{R}^p} R_\alpha(\xi, p_t(\xi)) p_t(\xi) d\xi \\ R_\beta(\xi, p_t(\xi)) &\stackrel{\text{def}}{=} -\alpha (\mathbb{E}_{P_{V,Z}} [Z\varphi(V; \xi)])^2 \\ &\quad + \mathbb{E}_{\tilde{\xi} \sim p_t} \left[\left(\mathbb{E}_{P_V} [\varphi(V; \xi)\varphi(V; \tilde{\xi})] \right)^2 \right], \end{aligned}$$

- The Energy functional is precisely the population MMD, *i.e.*, $E_\alpha(p_t(\xi)) = \text{MMD}_{\mu_t}^\alpha [P, Q]$, where $p_t(\xi) = \mu_t(\xi)/d\xi$.

Simulations on Synthetic Data-Set

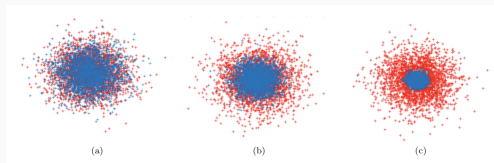


Figure 4: Panel (a): $\lambda = 0.1$, Panel (b): $\lambda = 0.5$, and Panel (c): $\lambda = 0.9$.

- Consider the problem of two sample test between two Gaussian distributions
 1. $P_0 = N(\mathbf{0}, (1 - \lambda)\mathbf{I}_{d \times d})$.
 2. $P_1 = N(\mathbf{0}, (1 + \lambda)\mathbf{I}_{d \times d})$.
- $\lambda \in [0, 1]$ controls the distance between these two distributions.
- Given samples $X_1, \dots, X_m \sim_{\text{i.i.d.}} P_0$ and $Z_1, \dots, Z_n \sim_{\text{i.i.d.}} P_1$, we want to decide between the following hypotheses
 1. **Null Hypothesis** $H_0: P_0 = P_1$ ($\lambda = 0$)
 2. **Alternative Hypothesis** $H_1: P_0 \neq P_1$ ($\lambda > 0$)

Simulations on Synthetic Data-Set

- We design a test statistics as below

$$\mathcal{T}(\{\mathbf{V}_i\}_{i=1}^m, \{\mathbf{W}_i\}_{i=1}^n) \stackrel{\text{def}}{=} \begin{cases} H_0 & \text{if } \widehat{\text{MMD}}_K[\{\mathbf{V}_i\}_{i=1}^m, \{\mathbf{W}_i\}_{i=1}^n] \leq \tau \\ H_1 & \text{if } \widehat{\text{MMD}}_K[\{\mathbf{V}_i\}_{i=1}^m, \{\mathbf{W}_i\}_{i=1}^n] > \tau, \end{cases}$$

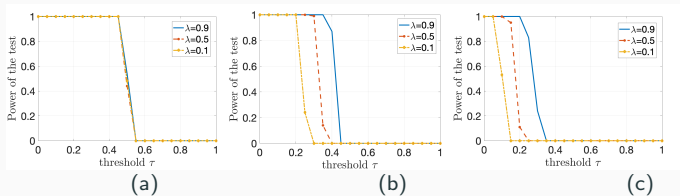


Figure 5: The statistical power, $\mathbb{P}(\text{reject } H_0 | H_1 \text{ is true})$, versus the threshold τ for the binary hypothesis testing via the unbiased estimator of the kernel MMD. Panel (a): Trained kernel using the two-phase procedure with the particle SGD and an auto-encoder, Panel (b): Trained kernel with an auto-encoder and a fixed Gaussian kernel with the bandwidth $\sigma = 1$, Panel (c): Untrained kernel without an auto-encoder.

Simulations on Real Data-Set: Qualitative Assessment

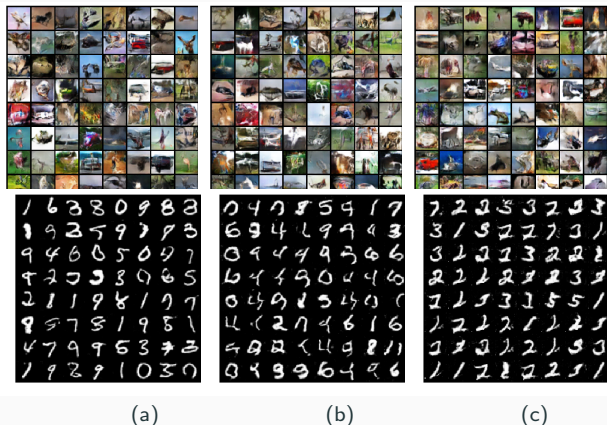


Figure 6: Sample generated images using CIFAR-10 and MNIST data-sets. Panel (a): Proposed MMD GAN with an *automatic* kernel selection via the particle SGD, Panel (b): MMD GAN with an auto-encoder optimization in conjunction with a mixed RBF Gaussian kernel, where the Gaussian bandwidths are *manually* tuned, Panel (c): MMD GAN with a single RBF Gaussian kernel with an auto-encoder optimization in conjunction with a single RBF Gaussian kernel where the Gaussian bandwidth is *manually* tuned.

Simulations on Real Data-Set: Quantitative Assessment

Method	FID (\downarrow)	IS (\uparrow)
MMD GAN (Gaussian)	67.244 ± 0.134	5.608 ± 0.051
MMD GAN (Mixture Gaussian)	67.129 ± 0.148	5.850 ± 0.055
Our Algorithm	65.059 ± 0.153	5.97 ± 0.046
Benchmark	-	11.237 ± 0.116

Table 1: Comparison of the quantitative performance measures of MMD GANs with different kernel learning approaches.

Conclusion

- Many machine learning tasks deal with the problem of learning a map between two domains.
- Machine learning systems divide into the supervised and unsupervised models based on the training samples. The hybrid version is often called semi-supervised model.
- Kernel methods provide an alternative method to deep learning to learn functions.
- However, there are model selection issues in kernel methods that need to be addressed. In this talk, we proposed a novel method to resolve those model selection issues.