

Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy

Dougal J. Sutherland^{1,2} Hsiao-Yu Tung² Heiko Strathmann¹ Soumyajit De¹ Aaditya Ramdas³ Alex Smola² Arthur Gretton¹

¹Gatsby Computational Neuroscience Unit, University College London

²School of Computer Science, Carnegie Mellon University

³Departments of EECS and Statistics, University of California at Berkeley

ICLR, 2017/ Presenter: Anant Kharkar

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

- Generative models produce data - attempt to match ground truth distribution
 - How can we distinguish real from generated data?
 - How to train generator network?
- Proposal: use Maximum Mean Discrepancy (MMD) to distinguish distributions

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

Divergences

- Comparison of two probability distributions
- Kullback-Leibler (KL)

$$D_{KL}(P||Q) = - \sum_i P(i) \log \frac{Q(i)}{P(i)}$$

- Jensen-Shannon

$$JSD(P||Q) = \frac{1}{2} KL(P||M) + \frac{1}{2} KL(Q||M)$$

$$M = \frac{1}{2}(P + Q)$$

- Integral probability metrics
 - Witness function distinguishes P from Q
 - Includes MMD

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - **Maximum Mean Discrepancy**
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

Maximum Mean Discrepancy

- Definition

$$MMD_k^2(P, Q) = \mathbb{E}_{x, x'}[k(x, x')] + \mathbb{E}_{y, y'}[k(y, y')] - 2\mathbb{E}_{x, y}[k(x, y)]$$

- Estimate

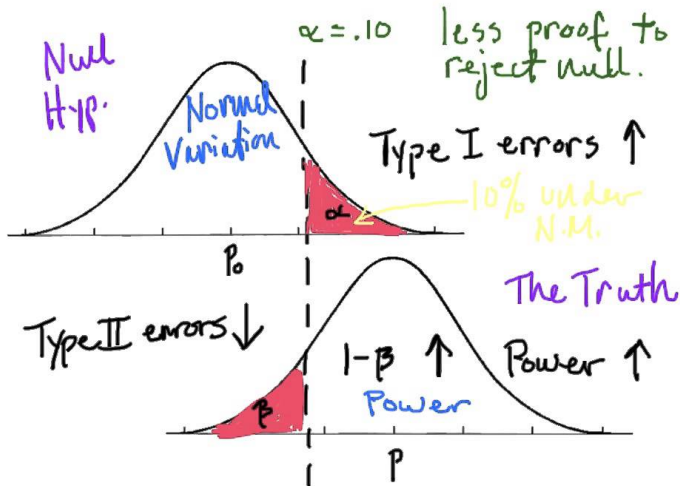
$$\widehat{MMD}_U^2(P, Q) = \frac{1}{\binom{m}{2}} \sum_{j \neq j'} k(X_j, X_{j'}) + \frac{1}{\binom{m}{2}} \sum_{i \neq i'} k(Y_i, Y_{i'}) - \frac{2}{\binom{m}{2}} \sum_{i \neq j} k(X_i, Y_j)$$

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - **Test Power**
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

Test Power

- Power: probability of rejecting H_0 given that H_A is true
- Measure of effectiveness of hypothesis



Hypothesis Test

- $H_0 : P = Q$
- $H_1 : P \neq Q$

$$\frac{\widehat{\text{MMD}}_{\text{U}}^2(X, Y) - \text{MMD}^2(P, Q)}{\sqrt{V_m(P, Q)}} \xrightarrow{D} \mathcal{N}(0, 1)$$

$$\begin{aligned} \Pr_1 \left(m \widehat{\text{MMD}}_{\text{U}}^2(X, Y) > \hat{c}_\alpha \right) &= \Pr_1 \left(\frac{\widehat{\text{MMD}}_{\text{U}}^2(X, Y) - \text{MMD}^2(P, Q)}{\sqrt{V_m(P, Q)}} > \frac{\hat{c}_\alpha/m - \text{MMD}^2(P, Q)}{\sqrt{V_m(P, Q)}} \right) \\ &\rightarrow \Phi \left(\frac{\text{MMD}^2(P, Q)}{\sqrt{V_m(P, Q)}} - \frac{c_\alpha}{m\sqrt{V_m(P, Q)}} \right) \end{aligned} \quad (4)$$

- $H_0 : P = Q, H_1 : P \neq Q$

$$t_k(P, Q) := \text{MMD}_k^2(P, Q) / \sqrt{V_m^{(k)}(P, Q)}$$

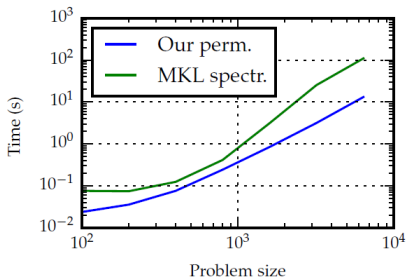
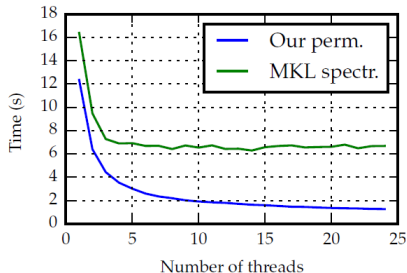
- $H_0 : P = Q, H_1 : P \neq Q$

$$t_k(P, Q) := \text{MMD}_k^2(P, Q) / \sqrt{V_m^{(k)}(P, Q)}$$

$$\begin{aligned} \hat{V}_m &:= \frac{2}{m^2(m-1)^2} \left(2\|\tilde{K}_{XX}e\|^2 - \|\tilde{K}_{XX}\|_F^2 + 2\|\tilde{K}_{YY}e\|^2 - \|\tilde{K}_{YY}\|_F^2 \right) \\ &\quad - \frac{4m-6}{m^3(m-1)^3} \left[\left(e^\top \tilde{K}_{XX}e \right)^2 + \left(e^\top \tilde{K}_{YY}e \right)^2 \right] + \frac{4(m-2)}{m^3(m-1)^2} (\|K_{XY}e\|^2 + \|K_{XY}^\top e\|^2) \\ &\quad - \frac{4(m-3)}{m^3(m-1)^2} \|K_{XY}\|_F^2 - \frac{8m-12}{m^5(m-1)} (e^\top K_{XY}e)^2 \quad (5) \\ &\quad + \frac{8}{m^3(m-1)} \left(\frac{1}{m} \left(e^\top \tilde{K}_{XX}e + e^\top \tilde{K}_{YY}e \right) (e^\top K_{XY}e) - e^\top \tilde{K}_{XX}K_{XY}e - e^\top \tilde{K}_{YY}K_{XY}^\top e \right). \end{aligned}$$

Optimization

- CPU cache optimization, multithreading
- Improved performance over Intel MKL spectral solver

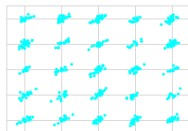
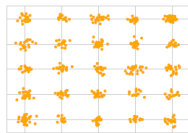


Outline

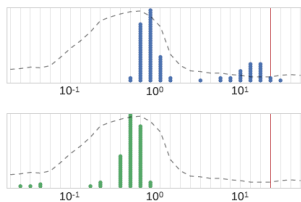
- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - Model Criticism
 - GAN
- 5 Summary

Synthetic Data

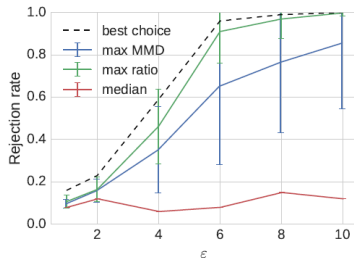
• Bandwidth selection for Gaussian RBF kernels - Blobs dataset



(a) Samples of size 500 with $\varepsilon = 6$ from P (top) and Q (bottom).



(b) Bandwidths chosen by maximizing $\widehat{\text{MMD}}_U^2$ (top, blue) and \hat{t} (bottom, green), as well as the median heuristic (red), for $\varepsilon = 6$. Gray lines show the power of each bandwidth: $\sigma = 0.67$ had power 96%, $\sigma = 10$ had 10%.



(c) Mean and standard deviation of rejection rate as ε increases. “Best choice” shows the mean power of the bandwidth with highest rejection rate for each problem.

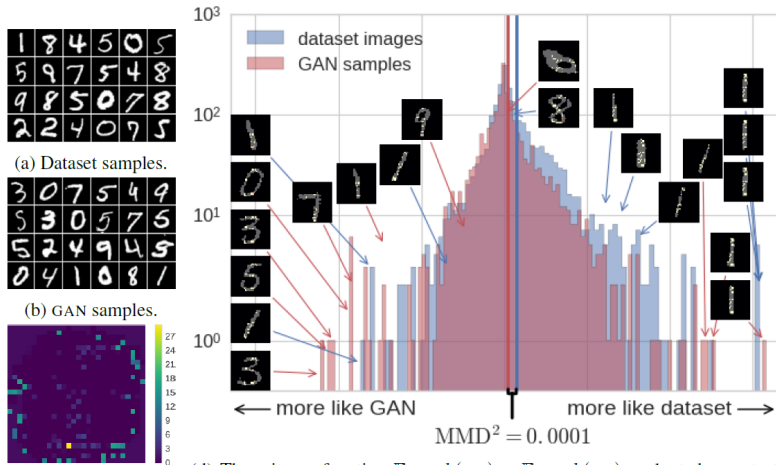
• Optimizing \hat{t} better than optimizing MMD

Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments
 - Synthetic Data
 - **Model Criticism**
 - GAN
- 5 Summary

Model Criticism

- Automatic relevance determination (ARD) kernel - MNIST



Outline

- 1 Introduction
 - Motivation
- 2 Background
 - Divergences
 - Maximum Mean Discrepancy
 - Test Power
- 3 Implementation
- 4 Experiments**
 - Synthetic Data
 - Model Criticism
 - GAN**
- 5 Summary

- Generative moment matching network (GMMN) uses MMD
- t -GMMN - minimizes t -statistic

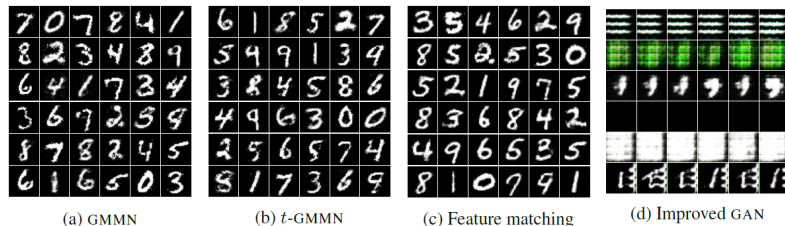


Figure 4: MNIST digits from various models. Part **d** shows six runs of the minibatch discrimination model of [Salimans et al. \(2016\)](#), trained without labels — the same model that, with labels, generated [Figure 3b](#). (The third row is the closest we got the model to generating digits without any labels.)

Summary

- MMD is a divergence metric
- Constructed MMD optimization t -test
- Propose MMD t -test as tool for GANs
 - Model criticism
 - GAN optimizer