Informative Features for Model Comparison

Heishiro Kanagawa² Wittawat Jitkrittum¹

¹Max Planck Institute for Intelligent Systems

When comparing complex generative models in high dimensions, the question to ask is not "which model is correct" (neither), or "which model is better," but rather "where does each model do better than the other?"

- **Given**: Two candidate models p, q, a sample $\{\mathbf{z}_i\}_{i=1}^n$ from an unknown distribution r.
- **Do**: Test H_0 : *p* fits sample better **vs** H_1 : *q* fits sample better. • **Propose**: Two new model comparison tests:
- 1. Rel-UME: Represent p, q by i.i.d. samples $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{y}_i\}_{i=1}^n$. 2. Rel-FSSD: p, q are unnormalized probability densities.
- Advantages:
- 1. Nonparametric: mild assumptions on p, q. Domain $\mathcal{X} \subseteq \mathbb{R}^d$.
- 2. Linear-time: $\mathcal{O}(n)$ runtime complexity. Fast. \bigcirc
- 3. Informative: show where q fits better than p (or vice versa) with a set of points (features).
- Test power matches that of the state-of-the-art quadratic-time relative MMD test [Bounliphone et al., 2015].

Test Statistics

• Let $D_V(p,r)$:= distance between p,r measured at V = $\{\mathbf{v}_1, \ldots, \mathbf{v}_J\}$ (features). D_V can be UME or FSSD. • Statistic $S := D_V(p, r) - D_V(q, r) \implies H_0: S \le 0 \text{ vs } H_1: S > 0.$ • $D_V(p, r)$ = average evaluation of the squared witness function $= \frac{1}{I} \sum_{i=1}^{J} \text{witness}_{p,r}^2(\mathbf{v}_j).$

UME (Unnormalized Mean Embeddings) [Jitkrittum et al., 2016] p / witness_{*p*,*r*}(**v**) = $\mathbb{E}_{\mathbf{x}\sim p}k(\mathbf{x}, \mathbf{v}) - \mathbb{E}_{\mathbf{z}\sim r}k(\mathbf{z}, \mathbf{v})$ for some kernel $k(\mathbf{x}, \mathbf{v})$.

FSSD (Finite-Set Stein Discrepancy) [Jitkrittum et al., 2017]

witness_{*p*,*r*}(**v**) = $\mathbb{E}_{\mathbf{x} \sim p} \boldsymbol{\xi}_{p}(\mathbf{x}, \mathbf{v})^{0} - \mathbb{E}_{\mathbf{z} \sim r} \boldsymbol{\xi}_{p}(\mathbf{z}, \mathbf{v})$ where $\xi_p(\mathbf{x}, \mathbf{v}) := \widetilde{k}(\mathbf{x}, \mathbf{v}) \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{v})$

(Normalizer of p not needed.)

Patsorn Sangkloy³

Relative Goodness-of-Fit Testing

Proposition. Empirical statistic S = S(V) for both Rel-UME and Rel-FSSD follows a normal distribution as $n \to \infty$.

- Estimate S. Reject H_0 if $S > \text{threshold} = (1 \alpha)$ -quantile of the normal. False rejection rate $< \alpha$ (asymptotically).
- Reject $H_0 \implies q$ is closer to r as measured at V.

Informative Features = V which maximizes

Test Power = $\mathbb{P}(\text{detect better fit of } q \mid q \text{ is better}).$

Equivalently, find V which maximizes the **power criterion**:

Power Criterion(V) := $\frac{S(V)}{\text{uncertainty}(V)} = \left(\frac{\text{signal}}{\text{noise}}\right)$,

- where uncertainty(V) = variance of S(V) under H_1 .
- $\mathcal{O}(n)$ complexity to evaluate power criterion. Fast.

Rel-UME and Rel-FSSD Power Criteria

Zero criterion $\implies p, q$ fit equally well. Extra mass of p = Missing mass of q. Rel-UME Rel-FSSD Criterion negative. *q* better here. *p* better here.

- Rel-UME: better model produces mass closer to the test sample from *r*.
- Rel-FSSD: better model has shape (given by $\nabla_{\mathbf{x}} \log p(\mathbf{x})$) and $\nabla_{\mathbf{v}} \log q(\mathbf{y})$) closer to r.

Contact: wittawat@tuebingen.mpg.de

James Hays³ Bernhard Schölkopf¹ Arthur Gretton²

²Gatsby Unit, University College London





• Set V = 40 (real) images of digit $i = 0, \ldots, 9$. • q is better at "1" and "5". p is slightly better at "3". Interpretable.



³Georgia Institute of Technology

criterion with n = 2000.