

Validation of generative models, especially SBI

- fundamental problem: if $x \sim p(x)$, there is no single correct outcome
 - traditional testing $\hat{X}_n = x_i^*$ does not work

⇒ must compare distributions, not individual outcomes

case (1) we have a test sample $\{x_i^* \sim p^*(x)\}_{i=1}^N$ or $\{x_i^* \sim p^*(x|y)\}_{i=1}^N$ for fixed y

easily generated by simulation $\hookrightarrow \int p^*(y) p^*(x|y) dy$

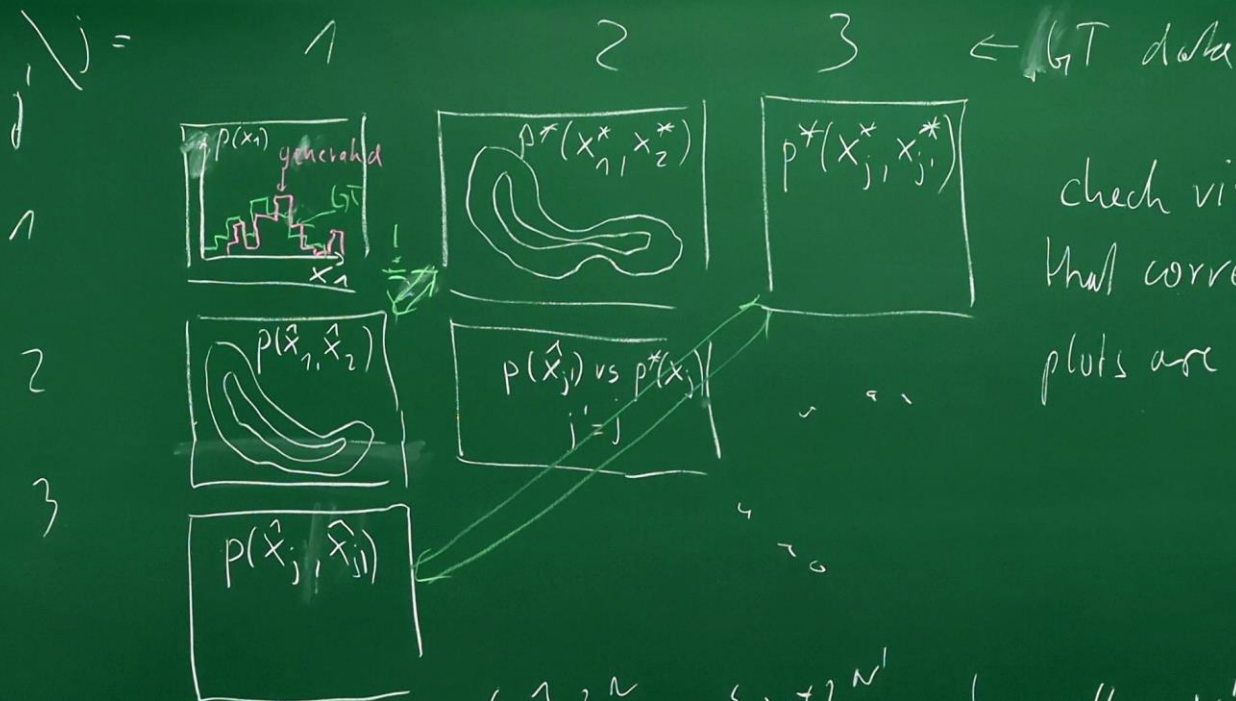
[more difficult for posterior $p(y|x) \approx p^*(y|x)$

⇒ to generate $\{y_i^* \sim p^*(y|x)\}_{i=1}^N$ for fixed x^{obs} , we need a classical algorithm like MCMC or ABC for $y \sim p^*(y|x^{\text{obs}})$]

- compare covariances of $\{\hat{X}_n\}$ and $\{X_n^*\}$: quick check, complete if p is Gauss means and
- plot marginal distributions in 1D and 2D for features $j, j' = 1, \dots, D$

example: $D=3$

synthetic data



various scores to compare the samples $\{x_i\}_{i=1}^N$ $\{x'_i\}_{i=1}^{N'}$ (usually, $N'=N$)

- MMD (maximum mean discrepancy)

$$MMD^2 = \frac{1}{N(N-1)} \sum_{i=1}^{N'} \sum_{i' \neq i}^N k(x_i^*, x_{i'}^*) + \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{i' \neq i}^N k(\hat{x}_i, \hat{x}_{i'})$$

$$- \frac{2}{N \cdot N'} \sum_{i=1}^{N'} \sum_{i'=1}^N k(x_i^*, \hat{x}_{i'})$$

with kernel $k(x, x')$

- Fréchet Inception Distance (FID) especially popular if X is an image
- idea: - compare means and covariance matrices via Fréchet distance
- if data are Gaussian with $\hat{\mu}, \hat{\Sigma}$ and μ^*, Σ^* for \hat{X} (Wasserstein-2 distance)

Fréchet distance can be computed analytically

$$FD = \|\hat{\mu} - \mu^*\|_2^2 + \text{tr} \left[\hat{\Sigma} + \Sigma^* - 2(\hat{\Sigma} \Sigma^*)^{1/2} \right]$$

for images: Gaussian assumption is not fulfilled $\hat{X}_i = h(\hat{X}_i), \hat{X}_i^* = h(X_i^*)$
 \Rightarrow instead calculate FD in some feature space
 fit $\hat{\mu}, \hat{\Sigma}$ and μ^*, Σ^* to \hat{X} and \hat{X}^* features

definition of $h(x)$: - typically, use a pre-trained network
 traditionally: Inception network for ImageNet

today: use some foundational model, e.g. CLIP (OpenAI)
 (trained with self-supervision on 10^8 images) DINOv2 (Meta)

1 famous algorithm for self-supervised learning: SIMCLR

- if no pre-trained model available: use a random (ly initialized) network

density & coverage (heuristic version of MMD²) [Nareem et al. 2020]

use nearest neighbor method, usually applied in a feature space to make
Euclidean distance plausible

define $B_k(x_i)$: ball with center x_i and radius $\|x_i - x_{i[k]}\|_2$

$\Rightarrow \mathbb{1}[\hat{x}' \in B_k(x_i)]$ is a kernel! k^{th} nearest neighbor of x_i

density = $\frac{1}{kN'} \sum_{i=1}^{N'} \sum_{i=1}^N \mathbb{1}[\hat{x}_i \in B_k(x_{i'}^*)]$ counts how many \hat{x}_i
each ball around $x_{i'}^*$ contains

\Rightarrow high when the \hat{x}_i are close to the $x_{i'}^*$

coverage = $\frac{1}{N'} \sum_{i=1}^{N'} \mathbb{1}[\exists i \text{ such that } \hat{x}_i \in B_k(x_{i'}^*)]$ counts how many
of the balls around $x_{i'}^*$

[recommended hyperparameters: $N=N'=10000$, $k=5$]

asymptotic limits when $p(x) = p^*(x)$ (distributions are identical)

$$\mathbb{E}[\text{density}] = 1$$

$$\mathbb{E}[\text{coverage}] = 1 - \frac{1}{2^k} \quad (if N \rightarrow \infty)$$

relation to MMD: $\text{density}(\hat{x}, \hat{x}) = \text{density}(x^*, x^*) = 1 \Rightarrow \text{MMD}^2 = 2 - \text{density}(\hat{x}, x^*) - \text{density}(x^*, \hat{x})$

case (2) we do not have a GT sample $\{X_i^* \}_{i=1}^N$

\Rightarrow must evaluate on the basis of $\{\hat{X}_i\}_{i=1}^N$ only, possible with a single instance $\hat{X} \sim p^*(X)$

[e.g. when learning the posterior, we have $\{(Y_i^* \sim p^*(Y), X_i^* = \phi(Y_i, \theta_1))\}_{i=1}^N$
then Y_i^* is a single GT example for $\hat{Y} \sim p(Y | X_i^*)$]

• check the diversity of the generated sample (without any GT)

Vendi score [Friedman & Dieng 2023]

kernel-based method like MMI diversity = $\{\hat{X}_i\}_{i=1}^N$ are all different, ideally cover entire $p^*(X)$

algorithm (1) calculate kernel (gram) matrix G :

$$G_{ii} = \frac{1}{N} k(\hat{X}_i, \hat{X}_i) \quad (\text{can also be done in a feature space } \mathcal{L}(X))$$

(1a) missing in the paper: centralize G as in kernel PCA (see MCE)

$$G \leftarrow C_N G C_N \quad C_N = \left(I_N - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right)$$

(2) calculate EV of G λ_i ($\text{tr } G = \sum_{i=1}^N \lambda_i = 1$)

(3) calculate the EV entropy $H(G) = -\sum_i \lambda_i \log \lambda_i$ (with $0 \cdot \log 0 = 0$)

(4) $VS = \exp(H(G)) \leq N$

VS acts as an effective rank of G

\Rightarrow properties: - if all points are far from each other: $G_{ii'} \approx 0$ for $i \neq i'$

$\Rightarrow VS = N$ (maximum diversity)

- if all points are equal: $G_{ii'} = \frac{1}{N}$ for all i, i'

$\Rightarrow VS = 1$ (minimum diversity)

choosing kernel bandwidth correctly is critical to avoid the extreme cases
[possibly, you can cheat with this]