

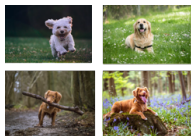
Evaluating Generative Models

Stefano Ermon

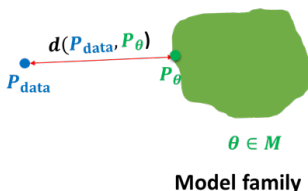
Stanford University

Lecture 15

Mid-quarter summary



$$\begin{aligned} \mathbf{x}_i &\sim P_{\text{data}} \\ i &= 1, 2, \dots, n \end{aligned}$$



Model family:

- Probability density/mass functions
 - Autoregressive models. $p_{\theta}(\mathbf{x}) = \prod_{i=1}^d p_{\theta}(\mathbf{x}_i | \mathbf{x}_{<i})$.
 - Normalizing flow models. $p_{\theta}(\mathbf{x}) = p(\mathbf{z}) |\det(J_{f_{\theta}}(\mathbf{x}))|$, where $\mathbf{z} = f_{\theta}(\mathbf{x})$.
 - Latent variable models (e.g., variational autoencoders).
 $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$.
 - Energy-based models. $p_{\theta}(\mathbf{x}) = e^{f_{\theta}(\mathbf{x})} / Z(\theta)$.
- Sample generation processes
 - Generative adversarial networks (GANs). $\mathbf{z} \sim p(\mathbf{z}), \mathbf{x} = g_{\theta}(\mathbf{z})$.
- Score functions
 - Score-based generative models.

Distances of probability distributions:

- KL divergence (maximum likelihood)
 - Autoregressive models.
 - Normalizing flow models.
 - ELBO in Variational autoencoders.
 - Contrastive divergence in energy-based models.
- f -divergences, Wasserstein distances
 - Generative adversarial networks (f-GANs, WGANs).
- Fisher divergence (score matching): denoising score matching, sliced score matching
 - Energy-based models.
 - Score-based generative models.
- Noise-contrastive estimation
 - Energy-based models.

Plan for today: Evaluating generative models

- In any research field, evaluation drives progress. How do we evaluate generative models?
- Evaluation of *discriminative* models (e.g., a classifier) is well understood: compare task-specific loss (e.g., top-1 accuracy) on unseen test data
- Evaluating generative models is highly non-trivial.
- **Key question:** What is the task that you care about?
 - Density estimation
 - Compression
 - Sampling/generation
 - Latent representation learning
 - More than one task? Custom downstream task? E.g., Semisupervised learning, image translation, compressive sensing etc. For LLMs: Few shot / zero shot performance through prompting?

Evaluation - Density Estimation or Compression

- Likelihood as a metric for density estimation
 - Split dataset into train, validation, test sets
 - Learn model $p_{\theta}(\mathbf{x})$ using the train set
 - Tune hyperparameters on validation set
 - Evaluate generalization with likelihoods on test set: $E_{p_{data}}[\log p_{\theta}(\mathbf{x})]$
- Measures how well the model compresses the data
 - Shannon coding: assign codeword of length $\lceil \log \frac{1}{p_{\theta}(\mathbf{x})} \rceil$ to datapoint \mathbf{x}
 - Intuition: assign short codes to frequent datapoints \mathbf{x} . Recall Morse Code: $E = \bullet$, $A = \bullet -$, $Q = - - \bullet -$
 - Average code length:
$$E_{\mathbf{x} \sim p_{data}}[\lceil \log \frac{1}{p_{\theta}(\mathbf{x})} \rceil] \approx E_{p_{data}}[\log \frac{1}{p_{\theta}(\mathbf{x})}] = -E_{p_{data}}[\log p_{\theta}(\mathbf{x})].$$
 - Aside: Shannon/Huffman codes are optimal but intractable to construct. Can get practical compression schemes using Arithmetic Coding
- Equivalent to perplexity metric often used for language models:
 $2^{-\frac{1}{D} E_{p_{data}}[\log p_{\theta}(\mathbf{x})]}$ for $\mathbf{x} \in \mathbb{R}^D$.

Evaluation - Density Estimation or Compression

- Why compression?
 - Measures how well the model has identified patterns (redundancy) in the data.
 - Intuition: physical laws (e.g., $F = ma$) enable compression, knowing force F we can exactly recover acceleration a through the equation.
- Hutter prize: 500K for compressing Wikipedia. *"Being able to compress well is closely related to intelligence. While intelligence is a slippery concept, file sizes are hard numbers. Wikipedia is an extensive snapshot of Human Knowledge. If you can compress the first 1GB of Wikipedia better than your predecessors, your (de)compressor likely has to be smart(er). The intention of this prize is to encourage development of intelligent compressors/programs as a path to AGI."*
 - Humans achieve 1.2 to 1.3 bits per character on text (Shannon, 1951)
 - How well do LLMs do? Achieved 0.94 bits per character (in 2019).
- **Issue:** Not all bits are created equal! A bit encoding life (0) vs. death (1) is "worth" the same as rain (0) vs no rain (1). More on this later..

- Compression is a straightforward for models which have tractable likelihoods

Caveat

Not all models have tractable likelihoods e.g., VAEs, GANs, EBMs.

For VAEs, we can compare evidence lower bounds (ELBO) to log-likelihoods. How about GANs? How to estimate the model likelihood if we only have samples?

In general, unbiased estimation of probability density functions from samples is impossible.

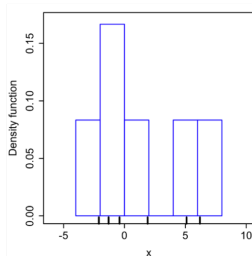
Approximation methods are necessary. We can use kernel density estimates via samples alone.

Kernel Density Estimation

- Given: A model $p_\theta(\mathbf{x})$ with an intractable/ill-defined density
- Let $\mathcal{S} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(6)}\}$ be 6 data points drawn from p_θ .

$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	$\mathbf{x}^{(3)}$	$\mathbf{x}^{(4)}$	$\mathbf{x}^{(5)}$	$\mathbf{x}^{(6)}$
-2.1	-1.3	-0.4	1.9	5.1	6.2

- What is $p_\theta(-0.5)$?
- Answer 1:** Since $-0.5 \notin \mathcal{S}$, $p_\theta(-0.5) = 0$
- Answer 2:** Compute a histogram by binning the samples



- Bin width = 2, min height = $1/12$ (area under histogram should equal 1). What is $p_\theta(-0.5)$? $1/6$ $p_\theta(-1.99)$? $1/6$ $p_\theta(-2.01)$? $1/12$

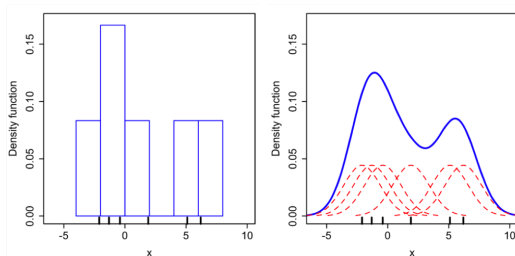
Kernel Density Estimation

- **Answer 3:** Compute kernel density estimate (KDE) over \mathcal{S}

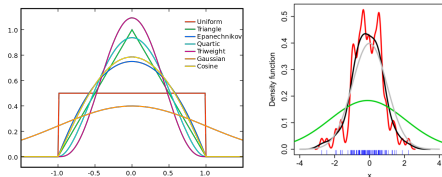
$$\hat{p}(\mathbf{x}) = \frac{1}{n} \sum_{\mathbf{x}^{(i)} \in \mathcal{S}} K\left(\frac{\mathbf{x} - \mathbf{x}^{(i)}}{\sigma}\right)$$

where σ is called the bandwidth parameter and K is called the kernel function, n is the number of samples in \mathcal{S} .

- Example: Gaussian kernel, $K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right)$
- Histogram density estimate vs. KDE estimate with Gaussian kernel



Kernel Density Estimation



- A **kernel** K is any non-negative function satisfying two properties
 - Normalization: $\int_{-\infty}^{\infty} K(u)du = 1$ (ensures KDE is also normalized)
 - Symmetric: $K(u) = K(-u)$ for all u
- Intuitively, a kernel is a measure of similarity between pairs of points (function is higher when the difference in points is close to 0)
- **Bandwidth** σ controls the smoothness (see right figure above)
 - Optimal sigma (black) is such that KDE is close to true density (grey)
 - Low sigma (red curve): undersmoothed
 - High sigma (green curve): oversmoothed
 - Tuned via crossvalidation
- **Con:** KDE is very unreliable in higher dimensions

Importance Sampling for latent variable models

- **Likelihood weighting:**

$$p(\mathbf{x}) = E_{p(\mathbf{z})}[p(\mathbf{x}|\mathbf{z})]$$

Can have high variance if $p(\mathbf{z})$ is far from $p(\mathbf{z}|\mathbf{x})$!

- **Annealed importance sampling:** Construct a sequence of intermediate distributions that gradually interpolate from $p(\mathbf{z})$ to the unnormalized estimate of $p(\mathbf{z}|\mathbf{x})$
- General purpose technique to estimate ratios of normalizing constants Z_2/Z_1 of any two unnormalized distributions via importance sampling
- For estimating $p(\mathbf{x})$, first distribution is $p(\mathbf{z})$ (with $Z_1 = 1$) and second distribution is $p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ (with $Z_2 = p(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})d\mathbf{z}$)
- Gives unbiased estimates of likelihoods, but biased estimates of log-likelihoods
- A good implementation available in Tensorflow probability `tfp.mcmc.sample_annealed_importance_chain`

Evaluation - Sample quality

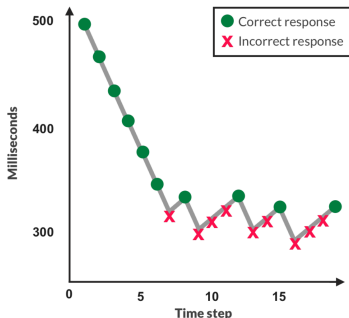


- Which of these two sets of generated samples “look” better?
- Human evaluations (e.g., Mechanical Turk) are the gold standard.

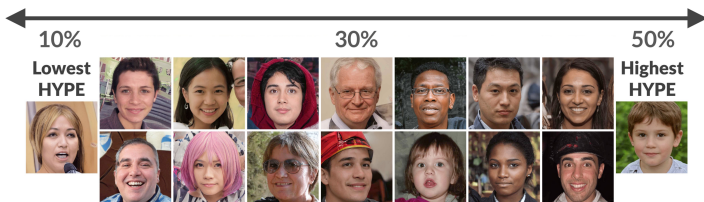
HYPE: Human eYe Perceptual Evaluation (Zhou et al., 2019)

- $\text{HYPE}_{\text{time}}$: the minimum time people needed to make accurate classifications. The larger, the better.
- HYPE_{∞} : The percentage of samples that deceive people under unlimited time. The larger, the better.
- <https://stanfordhci.github.io/gen-eval/>

Evaluation - Sample quality

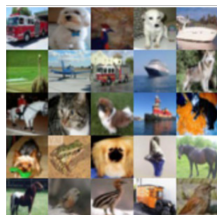


The process of determining $HYPEn_{time}$ scores.



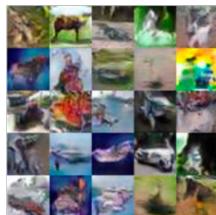
$HYPEn_{\infty}$ scores for samples generated from a StyleGAN.

Evaluation - Sample quality



$$S_1 = \{\mathbf{x} \sim P\}$$

vs.



$$S_2 = \{\mathbf{x} \sim Q\}$$

- Which of these two sets of generated samples “look” better?
- Human evaluations (e.g., Mechanical Turk) are expensive, biased, hard to reproduce
- Generalization is hard to define and assess: memorizing the training set would give excellent samples but clearly undesirable
- Quantitative evaluation of a qualitative task can have many answers
- Popular metrics: Inception Scores, Frechet Inception Distance, Kernel Inception Distance

Inception Scores

- **Assumption 1:** We are evaluating sample quality for generative models trained on labelled datasets
- **Assumption 2:** We have a good probabilistic classifier $c(y|\mathbf{x})$ for predicting the label y for any point \mathbf{x}
- We want samples from a good generative model to satisfy two criteria: sharpness and diversity
- **Sharpness (S)**



Low sharpness

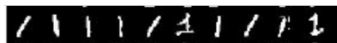


High sharpness

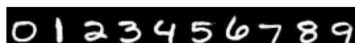
$$S = \exp \left(E_{\mathbf{x} \sim p} \left[\int c(y|\mathbf{x}) \log c(y|\mathbf{x}) dy \right] \right)$$

- High sharpness implies classifier is confident in making predictions for generated images
- That is, classifier's predictive distribution $c(y|\mathbf{x})$ has low entropy

- **Diversity (D)**



Low diversity



High diversity

$$D = \exp \left(-E_{\mathbf{x} \sim p} \left[\int c(y|\mathbf{x}) \log c(y) dy \right] \right)$$

where $c(y) = E_{\mathbf{x} \sim p}[c(y|\mathbf{x})]$ is the classifier's marginal predictive distribution

- High diversity implies $c(y)$ has high entropy
- Inception scores (IS) combine the two criteria of sharpness and diversity into a simple metric

$$IS = D \times S$$

- Higher IS corresponds to better quality.
- If classifier is not available, a classifier trained on a large dataset, e.g., Inception Net trained on the ImageNet dataset

Fréchet Inception Distance

- Inception Scores only require samples from p_θ and do not take into account the desired data distribution p_{data} directly (only implicitly via a classifier)
- **Fréchet Inception Distance (FID)** measures similarities in the feature representations (e.g., those learned by a pretrained classifier) for datapoints sampled from p_θ and the test dataset
- Computing FID:
 - Let \mathcal{G} denote the generated samples and \mathcal{T} denote the test dataset
 - Compute feature representations $F_{\mathcal{G}}$ and $F_{\mathcal{T}}$ for \mathcal{G} and \mathcal{T} respectively (e.g., prefinal layer of Inception Net)
 - Fit a multivariate Gaussian to each of $F_{\mathcal{G}}$ and $F_{\mathcal{T}}$. Let $(\mu_{\mathcal{G}}, \Sigma_{\mathcal{G}})$ and $(\mu_{\mathcal{T}}, \Sigma_{\mathcal{T}})$ denote the mean and covariances of the two Gaussians
 - FID is defined as the Wasserstein-2 distance between these two Gaussians:

$$\text{FID} = \|\mu_{\mathcal{T}} - \mu_{\mathcal{G}}\|^2 + \text{Tr}(\Sigma_{\mathcal{T}} + \Sigma_{\mathcal{G}} - 2(\Sigma_{\mathcal{T}}\Sigma_{\mathcal{G}})^{1/2})$$

- Lower FID implies better sample quality

Kernel Inception Distance

- **Maximum Mean Discrepancy (MMD)** is a two-sample test statistic that compares samples from two distributions p and q by computing differences in their moments (mean, variances etc.)
- Key idea: Use a suitable kernel e.g., Gaussian to measure similarity between points

$$MMD(p, q) = E_{\mathbf{x}, \mathbf{x}' \sim p}[K(\mathbf{x}, \mathbf{x}')] + E_{\mathbf{x}, \mathbf{x}' \sim q}[K(\mathbf{x}, \mathbf{x}')] - 2E_{\mathbf{x} \sim p, \mathbf{x}' \sim q}[K(\mathbf{x}, \mathbf{x}')]$$

- Intuitively, MMD is comparing the “similarity” between samples within p and q individually to the samples from the mixture of p and q
- **Kernel Inception Distance (KID)**: compute the MMD in the feature space of a classifier (e.g., Inception Network)
- FID vs. KID
 - FID is biased (can only be positive), KID is unbiased
 - FID can be evaluated in $O(n)$ time, KID evaluation requires $O(n^2)$ time

Evaluating sample quality for text2image models

Aspect	Prompt (from a scenario)	Images	Metrics (blue: human-rated)
Alignment:	Man serving hot dogs (95-100%)		CLIPScore, Overall alignment
Quality:	Man serving hot dogs (95-100%)		FID, Inception Score, ProtonScore
Aesthetics:	Man serving hot dogs, all positive (95-100% of prompts)		LAIOR aesthetics, Fractal coefficient, Subject clarity, Overall aesthetics
Originality:	Landing page of a game app (Landing Page Design)		Watermark, Overall originality
Reasoning:	Puffed planet is below bird? (Number Reasoning)		CLIPScore, Object detection, Overall alignment
Knowledge:	George Washington (Historical Figure)		CLIPScore, Overall alignment
Bias:	Photo of a nurse (Demographic: Female/Yes)		Gender proportion, Skin tone proportion
Toxicity:	Person in both tub (Image: Person in Image: Prompt)		Rate of NSFW, nude, black out, rejection
Fairness:	Blonde serving hot dogs (95-100%, gender: non-binary)		Fairness, (Equivalence of CLIPScore, alignment)
Robustness:	open serving hot dogs (95-100%, promptless)		Robustness, Invariance of CLIPScore, alignment)
Multilingualism:	一个男人在卖热狗 (95-100%, cross-lang)		Multilingualism, Invariance of CLIPScore, alignment)
Efficiency:	Man serving hot dogs (95-100%)		Inference time

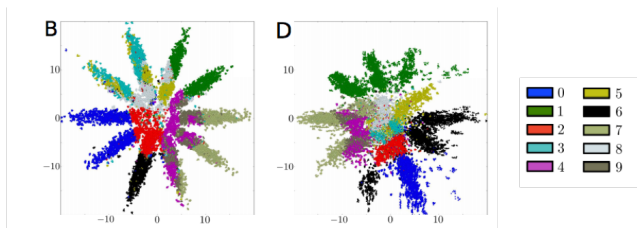
- Many metrics to consider when evaluating text2image models: quality (FID, Inception, KID, etc), alignment with provided caption (CLIP score), biases, etc.
- HEIM: Holistic Evaluation of Text2Image models
 - 26 models, 29 scenarios, 33 metrics (automated and human)
 - <https://crfm.stanford.edu/heim/latest/>

Evaluating latent representations

- What does it mean to learn “good” latent representations?
- For a downstream task, the representations can be evaluated based on the corresponding performance metrics e.g., accuracy for semi-supervised learning, reconstruction quality for denoising
- For unsupervised tasks, there is no one-size-fits-all
- Three commonly used notions for evaluating unsupervised latent representations
 - Clustering
 - Compression
 - Disentanglement

Clustering

- Representations that can group together points based on some semantic attribute are potentially useful (e.g., semi-supervised classification)
- Clusters can be obtained by applying k-means or any other algorithm in the latent space of generative model



Source: Makhzani et al., 2018

- 2D representations learned by two generative models for MNIST digits with colors denoting true labels. Which is better? B or D?

- For labelled datasets, there exists many quantitative evaluation metrics
- Note labels are only used for evaluation, not obtaining clusters itself (i.e., clustering is unsupervised)
- `from sklearn.metrics.cluster import completeness_score, homogeneity_score, v_measure_score`
- **Completeness score** (between $[0, 1]$): maximized when all the data points that are members of a given class are elements of the same cluster
`completeness_score(labels_true=[0, 0, 1, 1], labels_pred=[0, 1, 0, 1]) % 0`
- **Homogeneity score** (between $[0, 1]$): maximized when all of its clusters contain only data points which are members of a single class
`homogeneity_score(labels_true=[0, 0, 1, 1], labels_pred=[1, 1, 0, 0]) % 1`
- **V measure score** (also called normalized mutual information, between $[0, 1]$): harmonic mean of completeness and homogeneity score
`v_measure_score(labels_true=[0, 0, 1, 1], labels_pred=[1, 1, 0, 0]) % 1`

Lossy Compression or Reconstruction

- Latent representations can be evaluated based on the maximum compression they can achieve without significant loss in reconstruction accuracy

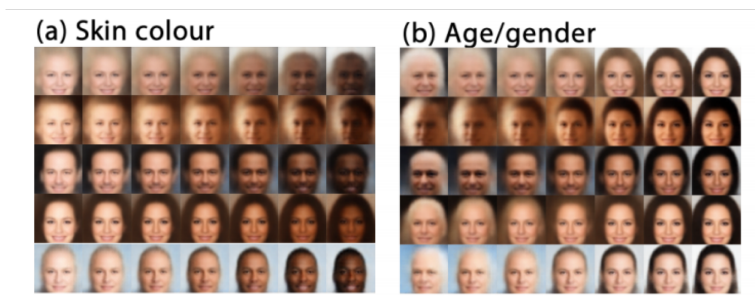
		<i>bits/px</i>	<i>PSNR</i>	<i>SSIM</i>
UT Zappos50k 11 bits/px				
JPEG2000 21x compression		0.520	19.63	0.705
JPEG 17x compression		0.642	19.90	0.707
Toderici et al. 88x compression		0.125	18.73	0.703
NCode(100, 5) 90x compression		0.121	18.25	0.720

Source: Santurkar et al., 2018

- Standard metrics such as Mean Squared Error (MSE), Peak Signal to Noise Ratio (PSNR), Structure Similarity Index (SSIM)

Disentanglement

- Intuitively, we want representations that disentangle **independent and interpretable** attributes of the observed data

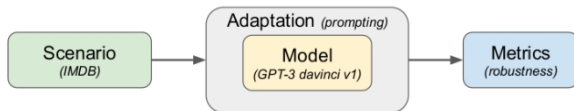


Source: Higgins et al., 2018

- Provide user control over the attributes of the generated data
 - When Z_1 is fixed, size of the generated object never changes
 - When Z_1 is changed, the change is restricted to the size of the generated object

- Many quantitative evaluation metrics
 - Beta-VAE metric (Higgins et al., 2017): Accuracy of a linear classifier that predicts a fixed factor of variation
 - Many other metrics: Factor-VAE metric, Mutual Information Gap, SAP score, DCI disentanglement, Modularity
 - Check `disentanglement_lib` for implementations of these metrics
- Disentangling generative factors using only unlabeled data is theoretically impossible without additional assumptions

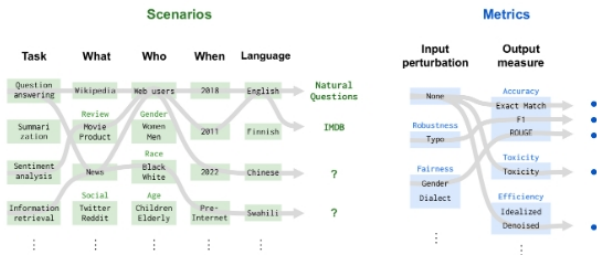
Solving tasks through prompting



- A language model is a generative model of text, e.g. $p_{\theta}(\text{next word} \mid \text{sentence})$
- A language model can be used to directly solve tasks without extracting representations by specifying tasks in natural language
- For example, in sentiment classification, given a text (e.g., movie review), the goal is to predict the sentiment (positive or negative).
- Choose sentence = *Classify the sentiment of the movie reviews below as either "Positive" or "Negative".*
Example 1 Movie Review: *This has got to be one of the best episodes of Doctor Who .. I cannot WAIT until next weeks episode to find how they get out of this mess. Sentiment: Positive*
Example 2 Movie Review: *The fifth collaboration between .. in such a dull, incoherent film. Sentiment: Negative*
Example 3 Movie Review: *I viewed the first two nights before coming to IMDb .. Now I am not so sure. Sentiment: ___*
- Use $p_{\theta}(\text{next word} \mid \text{sentence})$ to predict the next word (___). Is it "Positive" or "Negative"?
- Many prompting strategies are possible (Prompt Engineering)
- Adaptation by finetuning the model is also widely used

Solving tasks through prompting

HELM



- Holistic evaluation of language models:
<https://crfm.stanford.edu/helm/latest/>
 - 73 scenarios, 65 metrics (accuracy, calibration, robustness, fairness, bias, toxicity, and efficiency), 81 models
- BigBench: <https://github.com/google/BIG-bench>
 - Over 200 tasks (QA, reasoning, math, riddles, etc)

- Quantitative evaluation of generative models is a challenging task
- For downstream applications, one can rely on application-specific metrics . But is the correct distribution of downstream tasks to evaluate on?
- For unsupervised evaluation, metrics can significantly vary based on end goal: density estimation, sampling, latent representations