

GENERATIVE MODELS

FROM GAUSSIAN SAMPLING TO DIFFUSION MODELS

Alexandre Vérine,
Research Fellow, ENS-PSL
Université PSL

IASD
Université Dauphine-PSL

September 21, 2025

WHO AM I?

- ▶ Alexandre Vérine
- ▶ Research Fellow, ENS-PSL
- ▶ Université PSL
- ▶ Previously: PhD student at Dauphine (2021-2024)
- ▶ Research interests: generative models, evaluation methods, quality diversity trade-off
- ▶ All my slides and code are available on www.alexverine.com

ABOUT THIS COURSE

Outline of the course:

- ▶ Lecture 1 (22/09, 9h–12h15): From sampling to the first generative model
- ▶ Practical 1 (29/09, 9h–12h15): Building and training a VAE
- ▶ Lecture 2 (03/11, 9h–12h15): GANs (from the first GAN to f-GAN, WGAN, discriminator rejection sampling)
- ▶ Practical 2 (08/12, 9h–12h15): Training a GAN
- ▶ Lecture 3 (15/12, 9h–12h15): Diffusion Models (DDPM to EDM, score-based and classifier guidance)
- ▶ Practical 3 (26/01, 9h–12h15): Comparing ODE and SDE in small dimensions
- ▶ Lecture 4 (02/02, 9h–12h15): Evaluating generative models (IS, Precision and Recall) + Project presentation
- ▶ Last session (10/02, 9h–12h15): Student Presentations

ABOUT THIS COURSE

Project Presentation:

- ▶ Depending on the number of students, you will have to work in pairs or groups of 3.
- ▶ Same model architecture for all groups.
- ▶ Each group will have to choose a paper for training, regularizing, sampling the model.
- ▶ Each project will be tested every morning for a month.
- ▶ Final presentation on the last session (10/02/2024). Depending on the number of students, each group will have 5 to 7 minutes to present their work.
- ▶ Report to be handed the day before the project at 23:59.
- ▶ Report will be graded on the clarity of the presentation, the quality of the writing, the quality of the intuition behind experiments but not the results.
- ▶ Grade: 40% presentation, 60% report.

AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

- 1 What is sampling? 6**
 - 1.1 Definition 6
 - 1.2 Random and Pseudo-random sampling 18
 - 1.3 From uniform to Gaussian / Gaussian mixtures 22
- 2 Generative Models 38**
 - 2.1 Definition 38
 - 2.2 Implicit vs Explicit Models 38
 - 2.3 Divergences 43
 - 2.4 Examples of Generative Models 52
- 3 Types of Generative Models 54**
 - 3.1 Autoregressive Models 54
 - 3.2 Normalizing Flows 59
 - 3.3 Energy-Based Models 69
- 4 From an Autoencoder to a Generative Model 73**
 - 4.1 Autoencoder 73
 - 4.2 Variational Autoencoder 75
 - 4.3 Variants of VAE 80

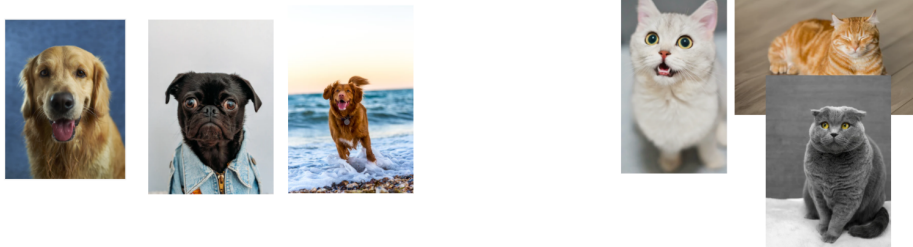
Part I

INTRODUCTION TO GENERATIVE MODELS

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :

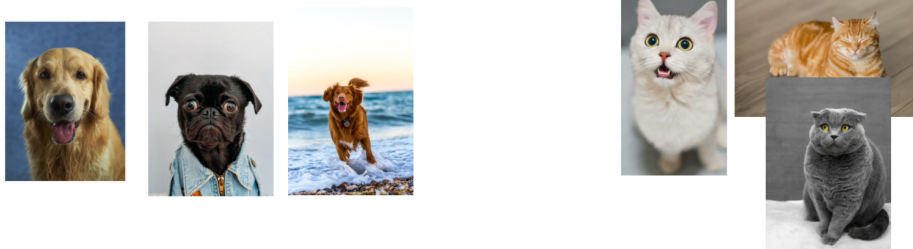


Data

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



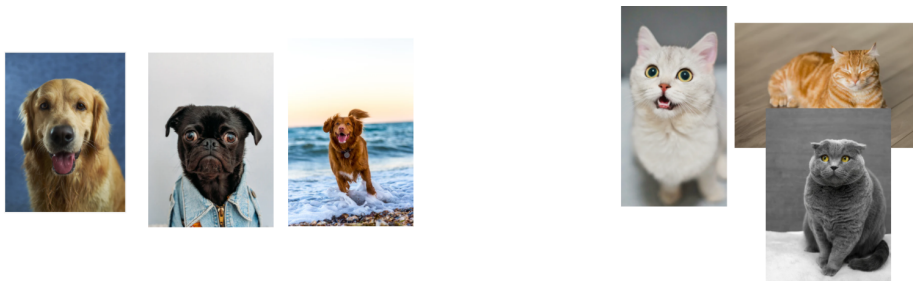
Data

What constraints should a generated image satisfy?

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

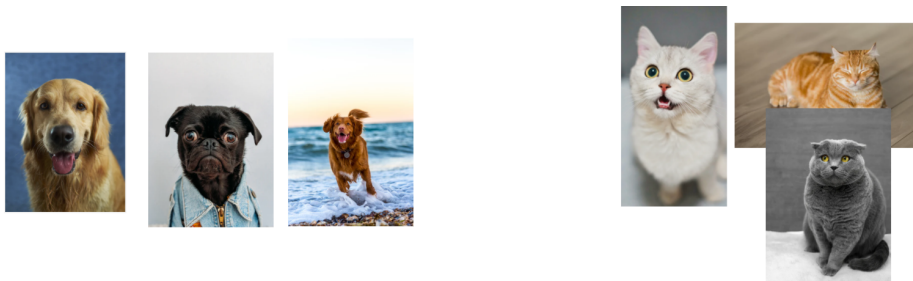


What constraints should a generated image satisfy?

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

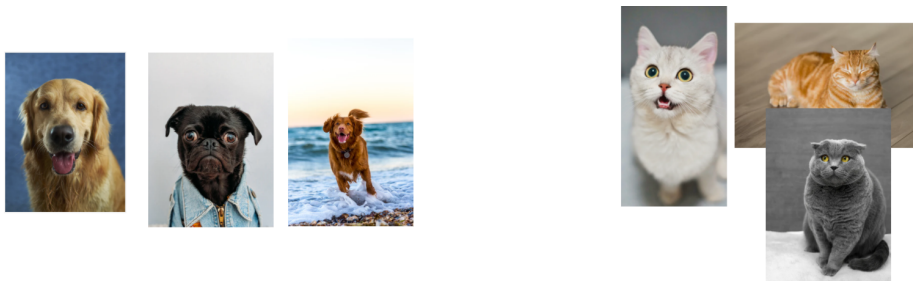


How do we ensure resemblance to real images?

WHAT IS SAMPLING?

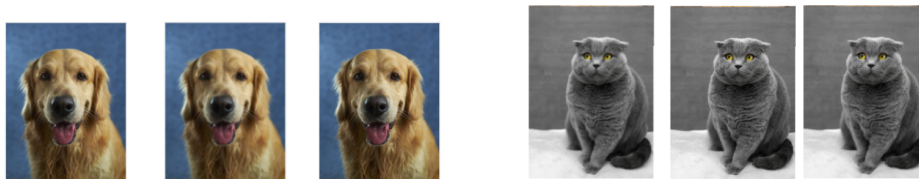
OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model

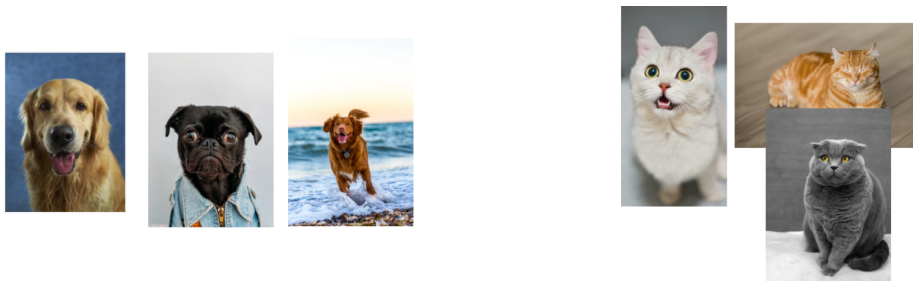


How do we ensure resemblance to real images?

WHAT IS SAMPLING?

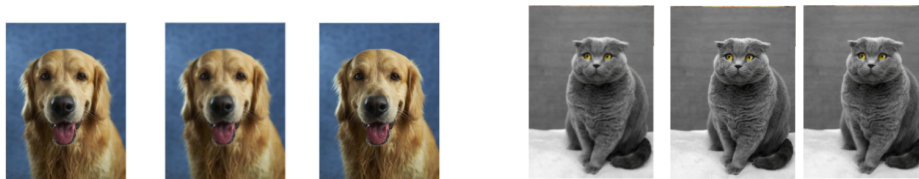
OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model

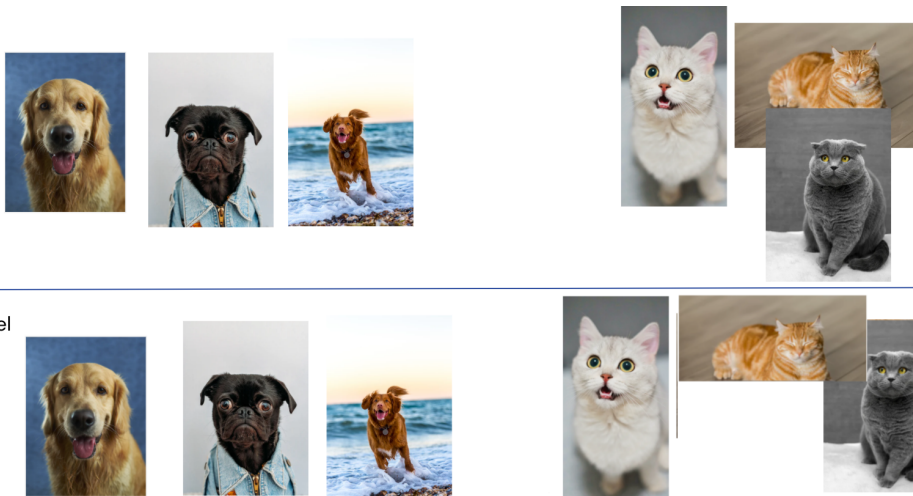


How to cover the diversity of the dataset?

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :

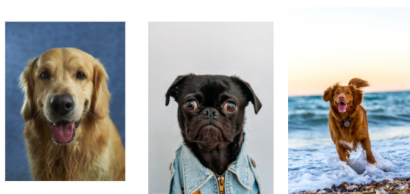


How to cover the diversity of the dataset?

WHAT IS SAMPLING?

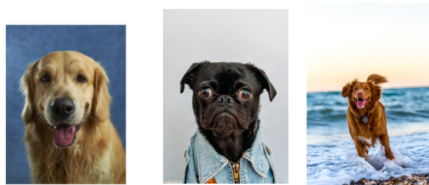
OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model

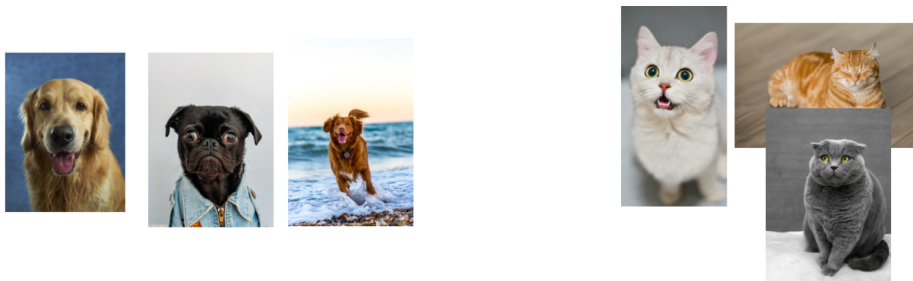


What about likelihood and novelty?

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model



What about likelihood and novelty?

WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model

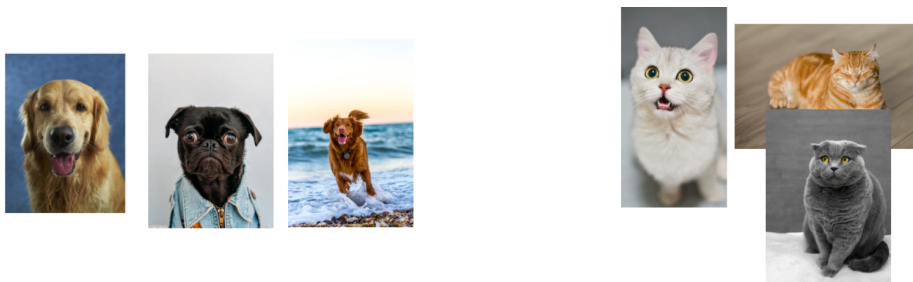


How to incorporate stochasticity in generation?

WHAT IS SAMPLING?

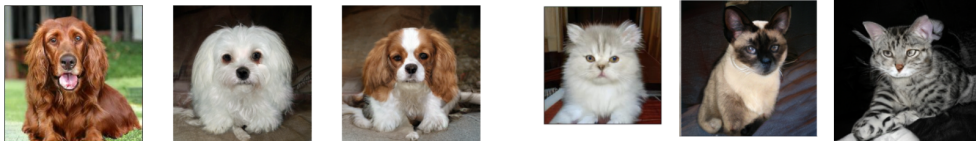
OBJECTIVE OF SAMPLING IMAGES

Images distributed under a given distribution P :



Data

Model



WHAT IS SAMPLING?

OBJECTIVE OF SAMPLING IMAGES

Constraints for sampling:

- ▶ *Resemblance*: Generated images should look like real ones.
- ▶ *Coverage*: Samples should represent the full data distribution.
- ▶ *Likelihood*: Samples should have high probability under the model.
- ▶ *Novelty*: Samples should not simply replicate training data.
- ▶ *Stochasticity*: Sampling should reflect inherent randomness.

WHAT IS SAMPLING?

WHY SAMPLING MATTERS

Definition:

Let P be a distribution on the sample space \mathcal{X} . The goal is to sample points under the distribution $p(x)$. Often, we consider conditional distributions $p(x|y)$ where y is some conditioning variable.

- ▶ Image generation : e.g., sampling new human faces (FFHQ, CelebA).
- ▶ Text generation : e.g., language models generating sentences.
- ▶ Image-to-image translation : e.g., translating day-to-night photos.
- ▶ Image-to-text generation : e.g., automatic image captioning.
- ▶ Text-to-image synthesis : e.g., generating images from text prompts.
- ▶ Text-to-text generation : e.g., machine translation.
- ▶ Speech-to-text transcription : e.g., transcribing audio to subtitles.
- ▶ Many other conditional generative tasks...

WHAT IS SAMPLING?

EXAMPLES: IMAGES-TO-IMAGES COLORIZATION

Example: Image Colorization

Original Image



Grayscale Image



Colorized Image



Task: Given a grayscale image, sample a plausible colored version.

WHAT IS SAMPLING?

EXAMPLES: IMAGES-TO-IMAGES INPAINTING

Example: Image Inpainting

Original Image



Masked Image



Filled Image



Task: Given an image with missing regions, sample a plausible completion.

WHAT IS SAMPLING?

EXAMPLES: IMAGES-TO-IMAGES UNCROPPING

Example: Image Uncropping

Original Image



Cropped Image



Filled Image



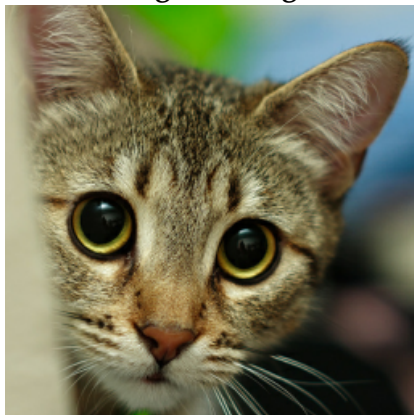
Task: Given an image with missing regions, sample a plausible completion.

WHAT IS SAMPLING?

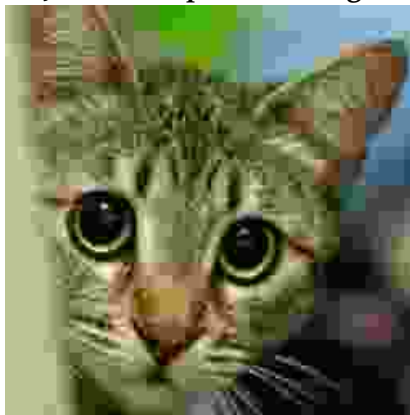
EXAMPLES: IMAGES-TO-IMAGES JPEG ARTIFACT REMOVAL

Example: Image JPEG Artifact Removal

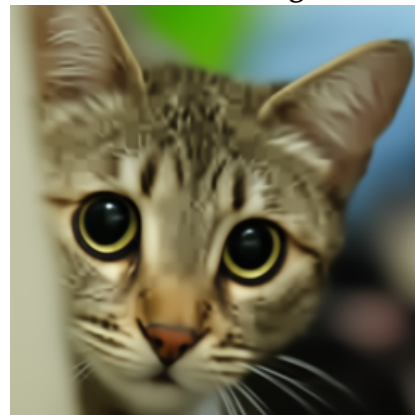
Original Image



JPEG Compressed Image



Restored Image



Task: Given a JPEG compressed image, sample a high-quality restoration.

WHAT IS SAMPLING?

EXAMPLES: TEXT-TO-IMAGE SYNTHESIS

Example: Text-to-Image Synthesis

Parti-350M



Parti-750M



Parti-3B



Parti-20B



A portrait photo of a kangaroo wearing an orange hoodie and blue sunglasses standing on the grass in front of the Sydney Opera House holding a sign on the chest that says Welcome Friends!

Task: Given a text description, sample a corresponding image.

WHAT IS SAMPLING?

EXAMPLES: TEXT-TO-TEXT TRANSLATION

Example: Text-to-Text Translation

Input (English):

A scenic view of a mountain during sunset.

A bustling city street at night.

A serene beach with palm trees.

Output (French):

Une vue pittoresque d'une montagne au coucher du soleil.

Une rue animée de la ville la nuit.

Une plage sereine avec des palmiers.

Task: Given an English text description, sample a French translation.

WHAT IS SAMPLING?

EXAMPLES: TEXT-TO-TEXT MATH PROBLEM SOLVING

Example: Text-to-Text Math Problem Solving

Input (Problem):

If a car travels at 60 mph for 2 hours, how far does it go?

What is the derivative of $x^2 + 3x + 2$?

Solve for x : $2x + 5 = 15$.

Output (Solution):

The car travels 120 miles.

The derivative is $2x + 3$.

The solution is $x = 5$.

Task: Given a text description, sample a corresponding solutions.

WHAT IS SAMPLING?

HOW DO WE ACTUALLY GENERATE AN IMAGE?

Prompt:

"A cute sloth holding a small treasure chest.
A bright golden glow is coming from the chest."

?



How should we generate something random corresponding to this?

WHAT IS SAMPLING?

HOW DO WE ACTUALLY GENERATE AN IMAGE?

Prompt:

"A cute sloth holding a small treasure chest.
A bright golden glow is coming from the chest."

?

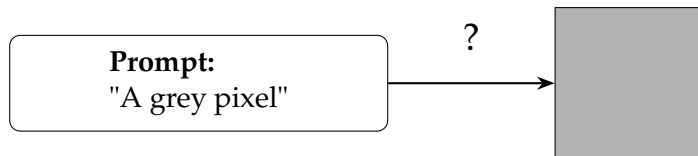


How should we generate something random corresponding to this?

It's very complex.

FROM COMPLEX TO SIMPLE

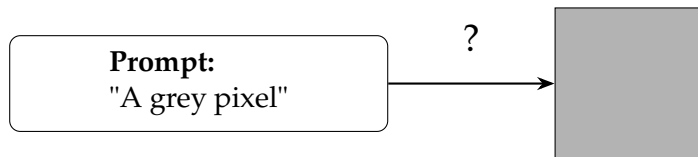
SIMPLIFYING THE PROBLEM



The simplest generation we can do: how can we do this?

FROM COMPLEX TO SIMPLE

SIMPLIFYING THE PROBLEM



The simplest generation we can do: how can we do this?

Even this is not trivial!

RANDOM VS PSEUDO-RANDOM SAMPLING

DEFINITIONS AND DIFFERENCES

Random (true randomness). Numbers generated from nondeterministic physical processes.

- ▶ Examples: quantum effects (photon arrival, electron tunneling), radioactive decay, thermal/shot noise.
- ▶ Properties: unpredictable, not reproducible from a finite state; entropy comes from physics.
- ▶ Uses: cryptography key generation, lotteries, high-stakes simulations.

Pseudo-random (PRNG). Numbers generated by deterministic algorithms.

- ▶ Defined by a recurrence/state update; *reproducible* given a seed.
- ▶ Aim for long period, good statistical tests, fast generation; not inherently cryptographically secure.
- ▶ Uses: ML training, Monte Carlo, graphics, games; seeding controls reproducibility.

RANDOM AND PSEUDO-RANDOM SAMPLING

SOURCES AND ALGORITHMS

True random sources (hardware/physical):

- ▶ Quantum RNGs (beam-splitter photon paths, vacuum fluctuations), radioactive decay counters.
- ▶ Electronic noise (thermal/Johnson–Nyquist, avalanche diode shot noise).
- ▶ External entropy (network jitter, disk timings) — lower quality, needs whitening.

Popular PRNG algorithms (software):

- ▶ Linear Congruential Generators (LCG), Xorshift / xoroshiro, PCG (permuted congruential generator).
- ▶ Mersenne Twister (MT19937): long period, equidistribution; standard in many libraries.
- ▶ Cryptographic PRNGs (ChaCha20-CTR, AES-CTR-DRBG) for security-sensitive use.

PSEUDO-RANDOM GENERATION

LINEAR CONGRUENTIAL GENERATOR (LCG)

Statement. We can generate *uniform* pseudo-random numbers efficiently with a simple PRNG (e.g., LCG) using a seed for reproducibility.

Algorithm 1: LCG (one-step)

Input: m, a, c, X_0

Output: $U_t \in (0, 1)$

$X_{t+1} \leftarrow (aX_t + c) \bmod m;$

$U_{t+1} \leftarrow X_{t+1}/m$

FROM UNIFORM TO MORE COMPLEX?

A NATURAL QUESTION

- ▶ OK, we can sample **uniform** variables quickly.
- ▶ **What about more complex distributions?** e.g., a Gaussian or a Gaussian mixture.
- ▶ Idea: *transform* uniforms into the target distribution. (Next: Box–Muller, Inverse CDF)

$$U \sim \text{Unif}(0, 1) \xRightarrow{\text{Box-Muller}} Z \sim \mathcal{N}(0, 1)$$

SAMPLING TRANSFORMATIONS

BOX-MULLER (UNIFORM \rightarrow GAUSSIAN)

Idea: map two i.i.d. uniforms (U_1, U_2) to two i.i.d. Gaussians (Z_1, Z_2).

- Polar transform: $R = \sqrt{-2 \ln U_1}$, $\Theta = 2\pi U_2$; then $Z_1 = R \cos \Theta$, $Z_2 = R \sin \Theta$.
- Produces *pairs* of normals; efficient reuse in vectorized code.
- Alternative: Marsaglia polar method avoids costly trig with rejection.

Algorithm 2: Box-Muller

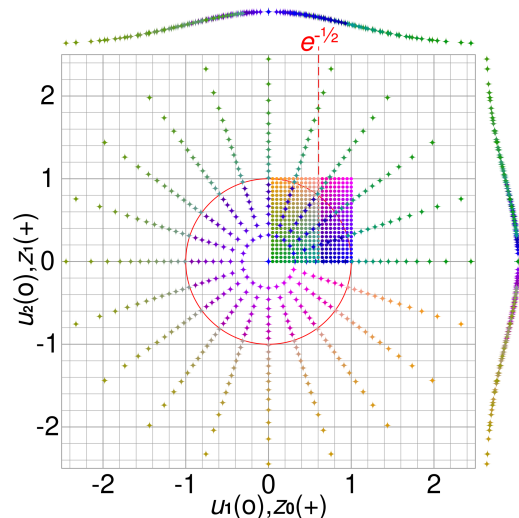
Output: $Z_1, Z_2 \sim \mathcal{N}(0, 1)$ i.i.d.

Draw $U_1, U_2 \sim \text{Unif}(0, 1)$ i.i.d.

$R \leftarrow \sqrt{-2 \log U_1}$

$\Theta \leftarrow 2\pi U_2$

$Z_1 \leftarrow R \cos \Theta$; $Z_2 \leftarrow R \sin \Theta$



INVERSE CDF

A.K.A. INVERSE TRANSFORM SAMPLING

- ▶ If $U \sim \text{Unif}(0, 1)$ and F is a CDF, then $X = F^{-1}(U)$ has CDF F .
- ▶ Requires F to be strictly increasing / invertible; otherwise use generalized inverse.
- ▶ In practice: precompute/discretize F , interpolate F^{-1} , or use spline/numerical root-finding.

INVERSE CDF

A.K.A. INVERSE TRANSFORM SAMPLING

- ▶ If $U \sim \text{Unif}(0, 1)$ and F is a CDF, then $X = F^{-1}(U)$ has CDF F .
- ▶ Requires F to be strictly increasing / invertible; otherwise use generalized inverse.
- ▶ In practice: precompute/discretize F , interpolate F^{-1} , or use spline/numerical root-finding.

Question: What is the issue with this method in high dimensions?

HIGH DIMENSIONAL SAMPLING

WHEN IS IT (RELATIVELY) EASY?

Fact: Sampling in high dimensions is generally hard, even if we know $p(x_1, \dots, x_N)$.

Special cases where it becomes easy(er):

- ▶ **Independence:** $p(x_1, \dots, x_N) = \prod_{i=1}^N p_i(x_i)$ — sample each coordinate independently.
- ▶ **Markov (autoregressive) chain:** $p(x_{1:N}) = p(x_1) \prod_{i=2}^N p(x_i | x_{i-1})$ — ancestral sampling along the chain.
- ▶ **Tree-structured Bayesian networks:** order variables topologically and sample parents \rightarrow children.

FROM UNIVARIATE TO MULTIVARIATE GAUSSIAN

DENSITIES AND PARAMETERS

Univariate (1D) Gaussian. For mean $\mu \in \mathbb{R}$ and standard deviation $\sigma > 0$,

$$X \sim \mathcal{N}(\mu, \sigma^2), \quad p(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

Parameters:

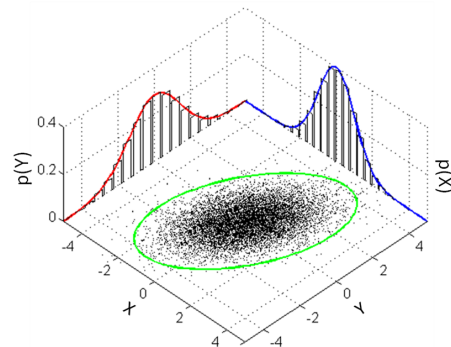
- ▶ μ (location): shifts the center of mass of the density.
- ▶ σ (scale): controls spread (variance σ^2).

Multivariate (d D) Gaussian. For mean $\mu \in \mathbb{R}^d$ and covariance $\Sigma \in \mathbb{R}^{d \times d}$ (symmetric PD),

$$X \sim \mathcal{N}(\mu, \Sigma), \quad p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right).$$

Parameters:

- ▶ μ (location vector): shifts the center in \mathbb{R}^d .
- ▶ Σ (covariance): encodes scale/shape/orientation; $\Sigma \succ 0$.



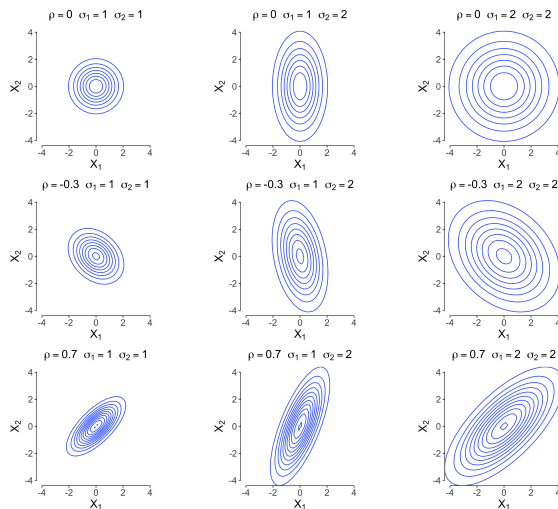
Elliptical contours illustrating Σ in \mathbb{R}^d .

FROM UNIVARIATE TO MULTIVARIATE GAUSSIAN

COVARIANCE TYPES

Special cases.

- Isotropic: $\Sigma = \sigma^2 \mathbf{I} \Rightarrow$ spherical contours.
- Diagonal: $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2) \Rightarrow$ independent coordinates.



Elliptical contours illustrating Σ in \mathbb{R}^d .

FROM $\mathcal{N}(0, \mathbf{I})$ TO $\mathcal{N}(\mu, \Sigma)$

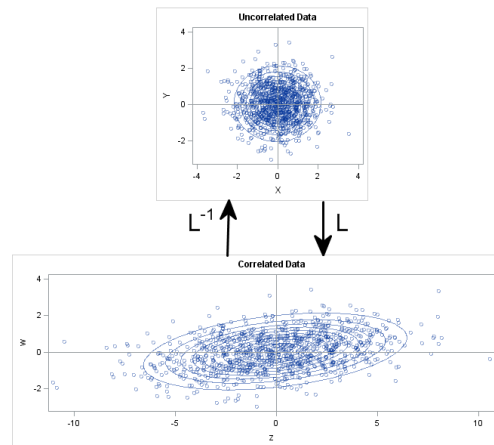
CHOLSKY FACTORIZATION AND SAMPLING

Cholesky factorization. For a symmetric positive-definite Σ , there exists a unique lower-triangular L with positive diagonal such that $\Sigma = LL^\top$.

Sampling recipe.

1. Draw $z \sim \mathcal{N}(0, \mathbf{I}_d)$ (independent standard normals).
2. Compute L such that $LL^\top = \Sigma$ (Cholesky).
3. Set $x \leftarrow \mu + Lz \Rightarrow$ then $x \sim \mathcal{N}(\mu, \Sigma)$.

Why it works. $\mathbb{E}[z] = 0$, $\text{Cov}(z) = \mathbf{I}$ and $\text{Cov}(Lz) = LIL^\top = \Sigma$.



Map unit sphere (standard normal) to ellipse via L , then translate by μ .

GAUSSIAN MIXTURES

DEFINITION AND INTUITION

Motivation. To increase the model complexity beyond a single Gaussian, we can *mix* several Gaussians; this yields a flexible, multi-modal density.

Definition. A Gaussian Mixture Model (GMM) with K components on \mathbb{R}^d is

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x \mid \mu_k, \Sigma_k), \quad \pi_k \geq 0, \quad \sum_k \pi_k = 1.$$

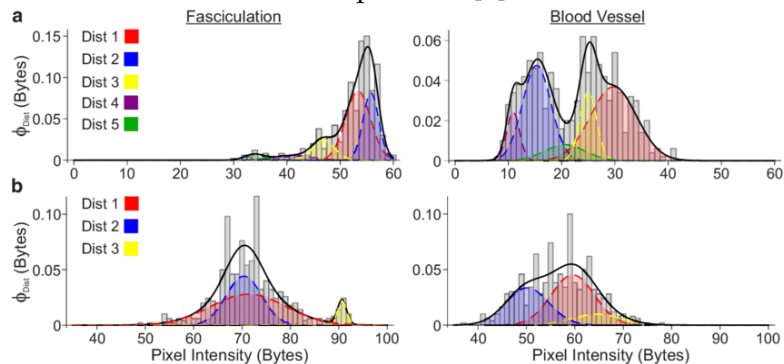
Parameters: weights $\{\pi_k\}$, means $\{\mu_k\}$, covariances $\{\Sigma_k\}$.

GAUSSIAN MIXTURES

WHY GAUSSIAN MIXTURES?

- *Expressivity*: universal approximator of smooth densities as K increases.

Example from [1]:



GAUSSIAN MIXTURES

SAMPLING FROM MIXTURE MODELS

Goal: draw $X \sim p(x) = \sum_k \pi_k \mathcal{N}(x \mid \mu_k, \Sigma_k)$.

Algorithm 3: Sampling from a GMM

Input: weights $\{\pi_k\}$, means $\{\mu_k\}$, covariances $\{\Sigma_k\}$

Output: $X \in \mathbb{R}^d$

Draw component index $K \sim \text{Categorical}(\pi_1, \dots, \pi_K)$

Draw $X \sim \mathcal{N}(\mu_K, \Sigma_K)$ (e.g., via Box–Muller/Cholesky)

Notes:

- ▶ Use prefix-sum table for the categorical draw; vectorize for batches.
- ▶ For Σ_K : diagonal for speed; Cholesky factor L with $LL^\top = \Sigma_K$ for full-cov.

TRAINING GAUSSIAN MIXTURES

MAXIMUM LIKELIHOOD — PRINCIPLE

Hypothesis: we observe i.i.d. data $x_1, \dots, x_N \sim P$ and posit a parametric model $p_\theta(x)$ (here, a GMM).

► **Population objective:**

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{X \sim P} [\log p_\theta(X)]$$

► **Empirical log-likelihood (practice):**

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_\theta(x_n) = \arg \max_{\theta} \log \prod_{n=1}^N p_\theta(x_n).$$

► **Why the log?** Turns products into sums; numerically stable.

Notes: LLN links empirical and population objectives; for GMMs, regularize Σ_k with a small $\epsilon \mathbf{I}$ to avoid degeneracy.

TRAINING GAUSSIAN MIXTURES

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

Objective (given data $\{x_n\}_{n=1}^N$):

$$\mathcal{L}(\Theta) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right) = \sum_{n=1}^N \underbrace{\text{LSE}_k \left(\log \pi_k + \log \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right)}_{\text{log-sum-exp}}.$$

where $\text{LSE}(z_1, \dots, z_K) := \log \sum_{k=1}^K e^{z_k}$.

- ▶ Constraints: $\pi \in \Delta^{K-1}$ (simplex), $\Sigma_k \succ 0$ (PD).
- ▶ Numerics: compute with the *log-sum-exp trick* (subtract \max_k) for stability.
- ▶ EM connection: LSE is a smooth max; E-step computes responsibilities as a softmax:

$$\gamma_{nk} = \text{softmax}_k \left(\log \pi_k + \log \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right).$$

TRAINING GAUSSIAN MIXTURES

GRADIENT DESCENT APPROACH

Parameterization for constraints:

- ▶ Weights: raw logits w_k with $\pi_k = \text{softmax}(w)_k$.
- ▶ Covariance: diagonal with $\sigma_{k,j}^2 = \exp(\alpha_{k,j})$; or full-cov via Cholesky $\Sigma_k = L_k L_k^\top$.

Stabilization:

- ▶ Compute $\log \sum_k$ via *log-sum-exp*; clip α ; add ϵI to Σ_k .
- ▶ Mini-batch SGD/Adam; early stopping; multiple restarts.

Pseudocode (one step):

Algorithm 4: GD step for GMM MLE

Input: minibatch $\{x_b\}$, current Θ

Compute $\ell_b = \log \sum_k \pi_k \mathcal{N}(x_b \mid \mu_k, \Sigma_k)$ for all b (log-sum-exp)

$J \leftarrow -\frac{1}{|B|} \sum_b \ell_b$

// negative log-likelihood

Backprop to get $\nabla_{\Theta} J$; update $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} J$ (Adam)

TRAINING GAUSSIAN MIXTURES

EM ALGORITHM — INTUITION

From LSE to EM. Using the LSE view, the softmax acts like a *soft argmax*: in the E-step we softly pick the near-maximum component; in the M-step we optimize parameters using those soft weights (points near the max matter more).

$$\gamma_{nk} = \text{softmax}_k \left(\log \pi_k + \log \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right).$$

Latent variables: introduce $z_{nk} \in \{0, 1\}$ (one-hot) with $\mathbb{P}(z_{nk} = 1) = \pi_k$.

- ▶ **E-step:** responsibilities $\gamma_{nk} = \mathbb{P}(z_{nk} = 1 \mid x_n, \Theta^{(t)})$ (soft assignments).
- ▶ **M-step:** maximize the expected complete-data log-likelihood

$$Q(\Theta, \Theta^{(t)}) = \sum_{n,k} \gamma_{nk} \left[\log \pi_k + \log \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right].$$

- ▶ **Guarantees:** each EM iteration non-decreasing in data log-likelihood; converges to a stationary point.

When to prefer EM vs GD: closed-form M-steps (fast, stable) vs flexible constraints/priors with GD.

TRAINING GAUSSIAN MIXTURES

EXPECTATION-MAXIMIZATION (EM) ALGORITHM (PSEUDO-CODE)

Algorithm 5: EM for Gaussian Mixture Models

Input: data $\{x_n\}_{n=1}^N$, number of components K

Output: parameters $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

Initialize π_k, μ_k, Σ_k (e.g., k-means)

repeat

 // E-step: responsibilities

for $n = 1$ **to** N **do**

for $k = 1$ **to** K **do**

$r_{nk} \leftarrow \pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k)$

end

$\gamma_{nk} \leftarrow r_{nk} / \sum_{j=1}^K r_{nj}$

end

 // M-step: parameter updates

$N_k \leftarrow \sum_{n=1}^N \gamma_{nk}$ **for** $k = 1, \dots, K$

$\pi_k \leftarrow N_k / N$

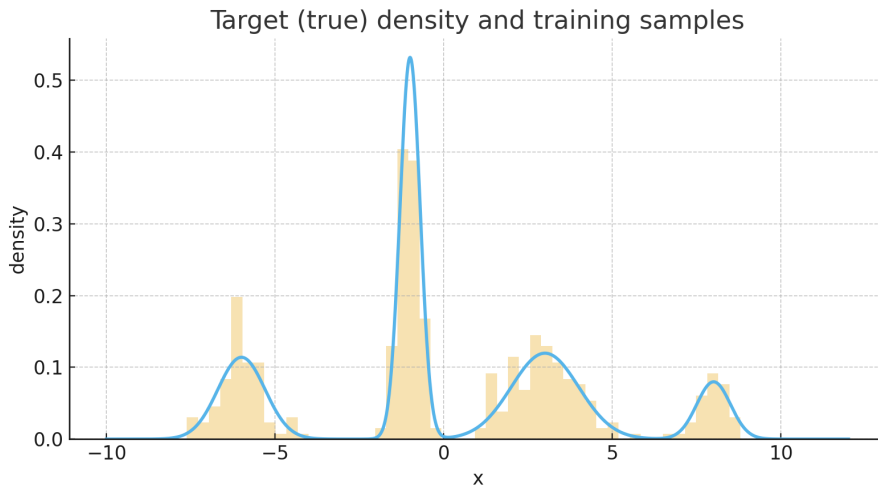
$\mu_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x_n$

$\Sigma_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top + \epsilon \mathbf{I}$

until convergence

GAUSSIAN MIXTURES

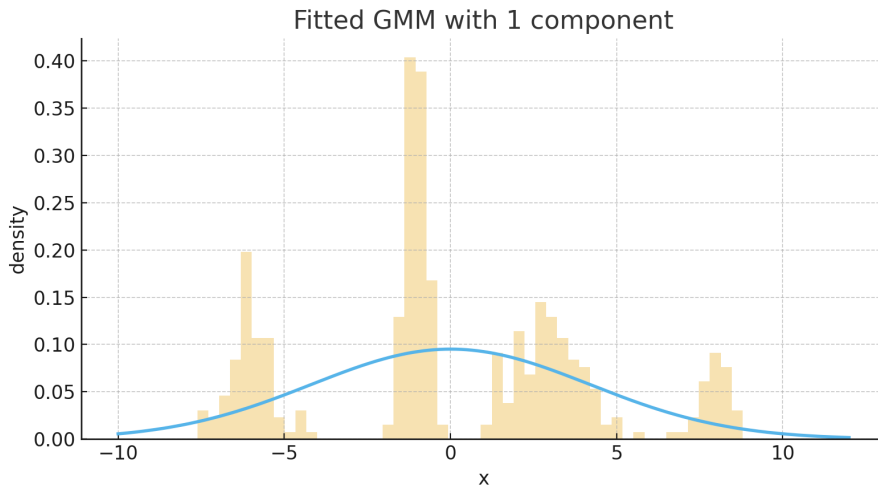
EXAMPLE



Fitting a GMM to 1D data points. In blue the density, in orange the histogram of data points.

GAUSSIAN MIXTURES

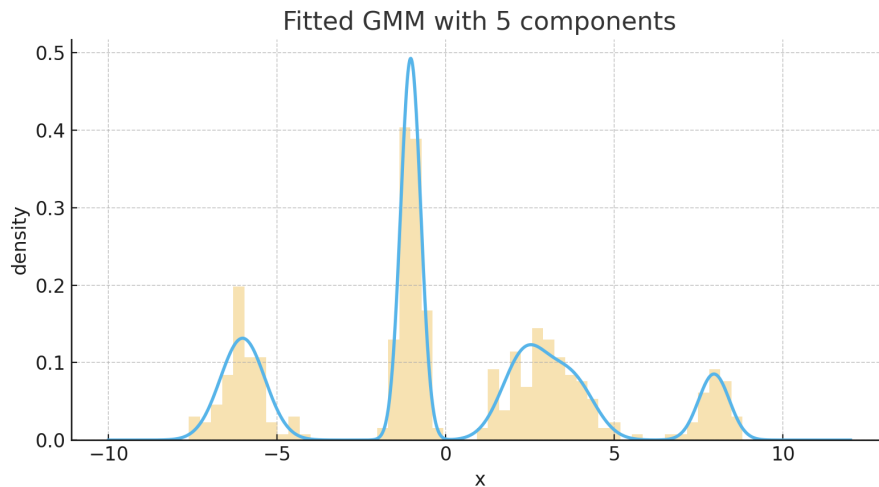
EXAMPLE



Fitting a GMM to 1D data points. $K = 1$.

GAUSSIAN MIXTURES

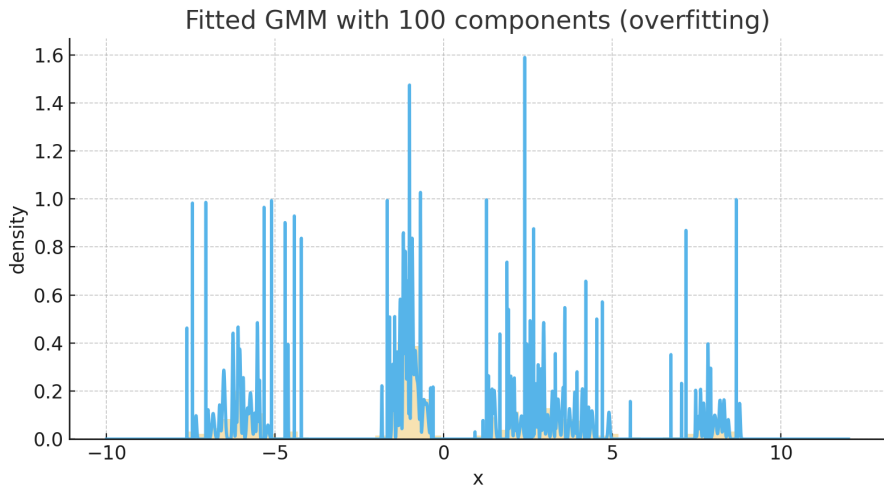
EXAMPLE



Fitting a GMM to 1D data points. $K = 5$.

GAUSSIAN MIXTURES

EXAMPLE



Fitting a GMM to 1D data points. $K = 150$.

GAUSSIAN MIXTURES

IN TERMS OF GENERATIVE MODELS

- ▶ *Resemblance*: If the modes are on the data, samples will look realistic.
- ▶ *Coverage*: More components \rightarrow better coverage of data distribution.
- ▶ *Likelihood*: GMMs provide explicit likelihoods for samples.
- ▶ *Novelty*: If not overfitting, samples can be novel.
- ▶ *Stochasticity*: Inherent randomness from component and Gaussian sampling.

GAUSSIAN MIXTURES

IN TERMS OF GENERATIVE MODELS

- ▶ *Resemblance*: If the modes are on the data, samples will look realistic.
- ▶ *Coverage*: More components → better coverage of data distribution.
- ▶ *Likelihood*: GMMs provide explicit likelihoods for samples.
- ▶ *Novelty*: If not overfitting, samples can be novel.
- ▶ *Stochasticity*: Inherent randomness from component and Gaussian sampling.

The GMM is a simple generative model! However, it has limitations in high dimensions and complex data (e.g., images).

GENERATIVE MODELS

TARGET DISTRIBUTION P AND MODEL \hat{P}_θ

Goal. Approximate the (unknown) data distribution P on \mathcal{X} with a *model family* $\{\hat{P}_\theta : \theta \in \Theta\}$.

Definition. A *generative model* is a probability distribution \hat{P}_θ over \mathcal{X} , parameterized by θ , together with a sampling procedure $x \sim \hat{P}_\theta$.

Types of generative models.

- ▶ **Explicit models:** provide a tractable density $\hat{p}_\theta(x)$ (e.g., GMMs, autoregressive, flows).
- ▶ **Implicit models:** define only a sampler $x = T_\theta(\varepsilon)$ with $\varepsilon \sim p(\varepsilon)$ (e.g., GANs, simulators).

IMPLICIT VS EXPLICIT MODELS

DEFINITIONS

Explicit (likelihood-based).

- ▶ Provide tractable density $\hat{p}_\theta(x)$ (or exact likelihood via change of variables).
- ▶ Examples: autoregressive (exact), flows (exact via Jacobian), GMMs (explicit), diffusion (via surrogate bounds).

What do explicit likelihood-based models allow us to do?

- ▶ *Out-of-distribution (OOD) detection*: Compute likelihoods for new samples; flag samples with low likelihood as OOD or anomalous.
- ▶ *Uncertainty quantification*: Assign probabilities to possible outcomes, enabling principled ways to measure uncertainty.
- ▶ *Anomaly detection*: Identify rare or unexpected events by their low likelihood under the model.
- ▶ *Model comparison*: Compare different generative models quantitatively using likelihoods or information criteria.
- ▶ *Principled training*: Enable maximum likelihood estimation and evaluation on held-out data.

IMPLICIT VS EXPLICIT MODELS

CHANGE OF VARIABLES FORMULA

Change of variables formula for densities:

Suppose $F : \mathcal{X} \rightarrow \mathcal{Z}$ is an *invertible* and *differentiable* function between open subsets of \mathbb{R}^d . Let $q(z)$ be a density on \mathcal{Z} . Then the induced density on \mathcal{X} is:

$$p(x) = q(F(x)) \cdot |\det \text{Jac}_F(x)|$$

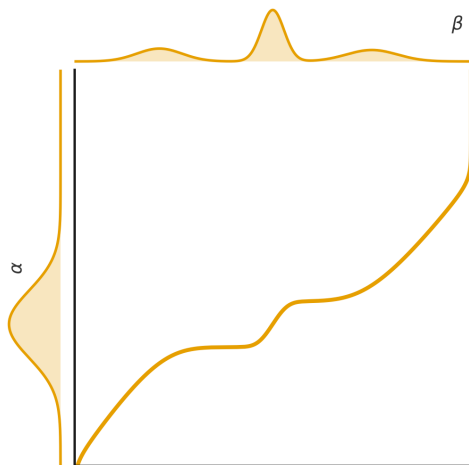
where $\text{Jac}_F(x)$ is the Jacobian matrix of F at x .

Conditions:

- ▶ F must be invertible (bijection) and differentiable, with differentiable inverse.
- ▶ The determinant of the Jacobian must be nonzero everywhere in the domain.

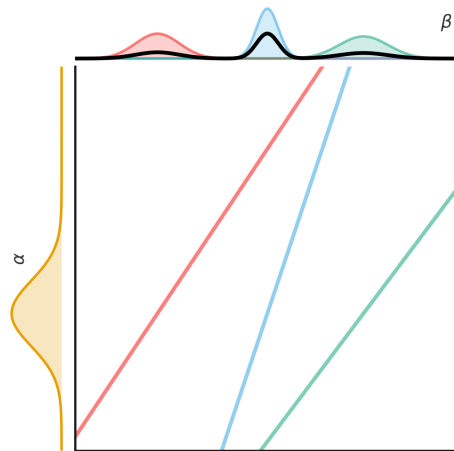
IMPLICIT VS EXPLICIT MODELS

CHANGE OF VARIABLES FORMULA



Deterministic Monge transport ($T = F\beta^{-1} \circ F\alpha$)

Illustration of the change-of-variables formula: mapping a simple distribution (e.g., gaussian) through a deterministic function F to obtain a more complex distribution.



Stochastic transport: three linear maps (left map in red)

Illustration of the change-of-variables formula: mapping a simple distribution (e.g., gaussian) through a stochastic function to obtain a more complex distribution.

IMPLICIT VS EXPLICIT MODELS

PUSH-FORWARD AND LIMITATIONS OF CHANGE OF VARIABLES

Invertibility and Differentiability: Not always possible

Most functions F from \mathcal{X} to \mathcal{Z} are *not* invertible and differentiable everywhere, especially when mapping between spaces of different dimensions or when F is not one-to-one.

- ▶ The change-of-variables formula applies only when F is a bijection between open subsets of \mathbb{R}^d and both F and F^{-1} are differentiable.
- ▶ In practice, many interesting mappings (e.g., neural networks with bottlenecks, dimensionality reduction) do not satisfy these conditions.

Push-forward distribution (formal definition)

Given a measurable function $F : \mathcal{Z} \rightarrow \mathcal{X}$ and a probability distribution Q on \mathcal{Z} , the **push-forward** distribution $P = F_{\#}Q$ on \mathcal{X} is defined by:

$$P(B) = Q(F^{-1}(B)) \quad \text{for any measurable set } B \subseteq \mathcal{X}.$$

That is, P is the distribution of $x = F(z)$ when $z \sim Q$.

Summary: The push-forward framework generalizes change-of-variables to cases where F may not be invertible or differentiable, but in those cases, we cannot write a simple density formula for $p(x)$.

DIVERGENCES

DEFINITION AND ESTIMATION TARGET

Divergence $\mathcal{D}(P\|Q)$. A non-negative functional with $\mathcal{D}(P\|Q) \geq 0$ and $\mathcal{D}(P\|Q) = 0$ iff $P = Q$. Not necessarily symmetric; no triangle inequality (*not* a distance).

Learning objective.

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{D}(P \parallel \hat{P}_\theta).$$

Examples next: KL, Total Variation, Wasserstein.

Goal: Find θ by minimizing a divergence $\mathcal{D}(P \parallel \hat{P}_\theta)$ or a surrogate (e.g., ELBO).

DIVERGENCES

KULLBACK–LEIBLER (FORWARD KL)

For $P \ll Q$ with densities p, q ,

$$\mathcal{D}_{\text{KL}}(P\|Q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E}_{X \sim P}[\log p(X) - \log q(X)].$$

- ▶ Asymmetric; mode-covering when used as $\mathcal{D}_{\text{KL}}(P\|Q)$ in many settings.
- ▶ **MLE link:** $\mathbb{E}_P[\log q_\theta(X)] = -\mathcal{D}_{\text{KL}}(P\|Q_\theta) - H(P)$.
- ▶ Requires $q > 0$ wherever $p > 0$ (absolute continuity).

DIVERGENCES

TOTAL VARIATION (TV)

$$\text{TV}(P, Q) = \sup_{A \subseteq \mathcal{X}} |P(A) - Q(A)| = \frac{1}{2} \int |p(x) - q(x)| dx.$$

- ▶ Metric on probability measures; bounded in $[0, 1]$.
- ▶ Interpretable as maximum test error gap over events.
- ▶ Hard to estimate directly in high dimension.

DIVERGENCES

WASSERSTEIN (OPTIMAL TRANSPORT)

For cost $c(x, y) = \|x - y\|$ and couplings $\Pi(P, Q)$,

$$\mathcal{W}_1(P, Q) = \inf_{\pi \in \Pi(P, Q)} \mathbb{E}_{(X, Y) \sim \pi} [\|X - Y\|].$$

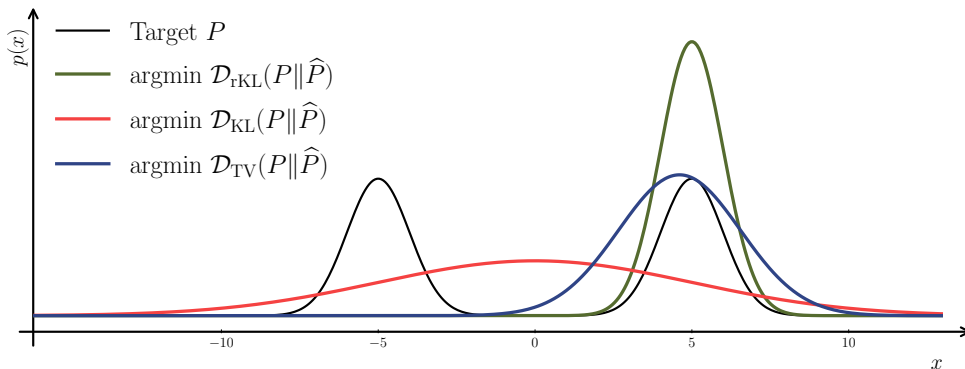
- ▶ Sensitive to the geometry of \mathcal{X} ; finite even with disjoint supports.
- ▶ Dual (Kantorovich–Rubinstein): $\mathcal{W}_1 = \sup_{\|f\|_{\text{Lip}} \leq 1} (\mathbb{E}_P[f] - \mathbb{E}_Q[f])$.
- ▶ Basis for WGAN objectives.

DIVERGENCES

KL VS REVERSE KL

Forward KL $\mathcal{D}_{\text{KL}}(P\|Q)$ tends to be *mode-covering*; **Reverse KL** $\mathcal{D}_{\text{KL}}(Q\|P)$ tends to be *mode-seeking*.

- ▶ Forward KL penalizes missing mass where $p > 0$ and $q \approx 0$ (heavy penalty).
- ▶ Reverse KL penalizes placing mass where $q > 0$ but $p \approx 0$; may ignore small modes.
- ▶ Choice impacts behavior of trained models.



HOW TO COMPUTE A DIVERGENCE?

KL AND MAXIMUM LIKELIHOOD

From KL to log-likelihood (derivation).

$$\begin{aligned}\mathcal{D}_{\text{KL}}(P\|Q_\theta) &= \int p(x) \log \frac{p(x)}{q_\theta(x)} dx \\ &= \mathbb{E}_{X \sim P} [\log p(X) - \log q_\theta(X)] \\ &= \underbrace{\mathbb{E}_P[\log p(X)]}_{-H(P)} - \mathbb{E}_P[\log q_\theta(X)] \\ &= -\mathbb{E}_P[\log q_\theta(X)] - H(P).\end{aligned}$$

Therefore,

$$\mathbb{E}_P[\log q_\theta(X)] = -\mathcal{D}_{\text{KL}}(P\|Q_\theta) - H(P),$$

and since $H(P)$ does not depend on θ ,

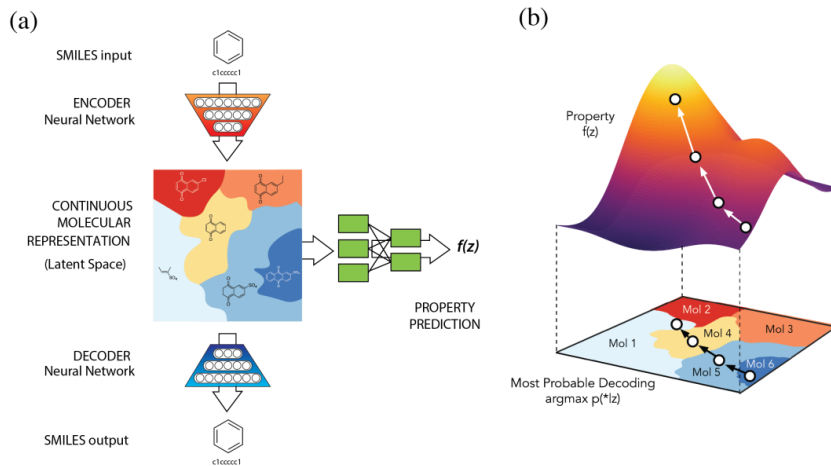
$$\theta^* = \arg \max_{\theta} \mathbb{E}_P[\log q_\theta(X)] \equiv \arg \min_{\theta} \mathcal{D}_{\text{KL}}(P\|Q_\theta).$$

HOW TO COMPUTE A DIVERGENCE?

LATENT VARIABLES AND LOW-DIMENSIONAL STRUCTURE

Idea. Introduce a *latent space* $\mathcal{Z} = \mathbb{R}^k$ with $k \ll d$ (data in \mathbb{R}^d) to capture low-dimensional structure.

- ▶ **Generative story:** sample $z \sim q(z)$ (prior), then $x \sim p_\theta(x | z)$ (decoder/model).
- ▶ **Inference:** approximate the posterior with $q_\phi(z | x)$ (encoder) since $p_\theta(z | x)$ is intractable.
- ▶ Benefits: compression, structure, disentanglement, controllable generation.



HOW TO COMPUTE A DIVERGENCE?

DERIVING THE ELBO

Goal: maximize $\log p_\theta(x)$ where $p_\theta(x) = \int p_\theta(x, z) dz$.
For any distribution $q_\phi(z | x)$,

$$\begin{aligned}\log p_\theta(x) &= \log \int q_\phi(z | x) \frac{p_\theta(x, z)}{q_\phi(z | x)} dz \\ &\geq \int q_\phi(z | x) \log \frac{p_\theta(x, z)}{q_\phi(z | x)} dz \\ &= \mathbb{E}_{q_\phi(z|x)} \left[\log p_\theta(x, z) - \log q_\phi(z | x) \right] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z) + \log p(z) - \log q_\phi(z | x)] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \mathcal{D}_{\text{KL}}(q_\phi(z | x) \| p(z))\end{aligned}$$

Define the **ELBO**:

$$\text{ELBO}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \mathcal{D}_{\text{KL}}(q_\phi(z | x) \| p(z)).$$

HOW TO COMPUTE A DIVERGENCE?

APPROXIMATE MLE AND ELBO

ELBO in KL form.

$$\text{ELBO}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)] - \mathcal{D}_{\text{KL}}(q_\phi(z | x) \| p(z)).$$

And the exact evidence decomposes as

$$\log p_\theta(x) \geq \text{ELBO}(\theta, \phi; x) - \mathcal{D}_{\text{KL}}(q_\phi(z | x) \| p_\theta(z | x)),$$

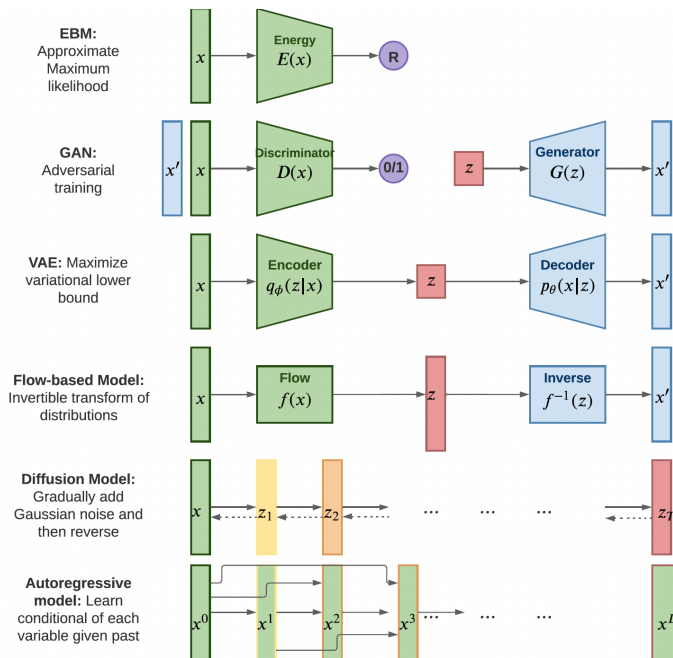
so maximizing ELBO *minimizes* the posterior KL.

What to optimize:

- ▶ **w.r.t. θ (decoder/model):** increase the reconstruction term $\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x | z)]$; use Monte Carlo gradients with reparameterized z .
- ▶ **w.r.t. ϕ (encoder/inference):** tighten the bound by driving $q_\phi(z | x)$ toward $p_\theta(z | x)$ (reducing the posterior KL above), while also respecting the regularizer $\mathcal{D}_{\text{KL}}(q_\phi(z | x) \| p(z))$.
- ▶ **Training recipe:** joint SGD on (θ, ϕ) with mini-batches and the reparameterization trick $z = \mu_\phi(x) + \sigma_\phi(x) \odot \varepsilon, \varepsilon \sim \mathcal{N}(0, \mathbf{I})$.

FAMILIES OF GENERATIVE MODELS

VISUAL TAXONOMY



MODEL LANDSCAPE

AT A GLANCE

Model	Density	Sampling	Training	Latents	Architecture	Discussed
ARM	Exact, fast	Slow	MLE	None	Sequential	Here
Flows	Exact, slow / fast	Slow	MLE	\mathbb{R}^d	Invertible	Here
EBM	Approx, slow	Slow	MLE-A	Optional	Discriminative	Here
VAE	LB, fast	Fast	MLE-LB	\mathbb{R}^m	Encoder–Decoder	Here and TP1
GAN	Jensen Approx	Fast	Min–max	\mathbb{R}^m	Generator–Discriminator	Session 2 and TP2
Diffusion	LB	Slow	MLE-LB	\mathbb{R}^d	Encoder–Decoder	Session 3 and TP3

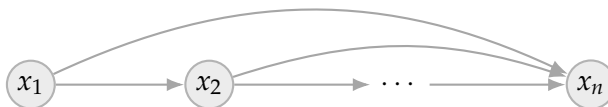
AUTOREGRESSIVE MODELS

LIKELIHOOD DECOMPOSITION

Chain rule factorization. For $x = (x_1, \dots, x_d)$ and any fixed ordering,

$$p_{\theta}(x) = \prod_{i=1}^d p_{\theta}(x_i \mid x_{<i}) \iff \log p_{\theta}(x) = \sum_{i=1}^d \log p_{\theta}(x_i \mid x_{<i}).$$

- ▶ The ordering (sequence order, raster scan for images, etc.) defines the conditional structure.
- ▶ Each factor is a *simple* conditional model (e.g., categorical over tokens/pixels, Gaussian for reals).
- ▶ Tractable likelihood: evaluation and gradients are exact.



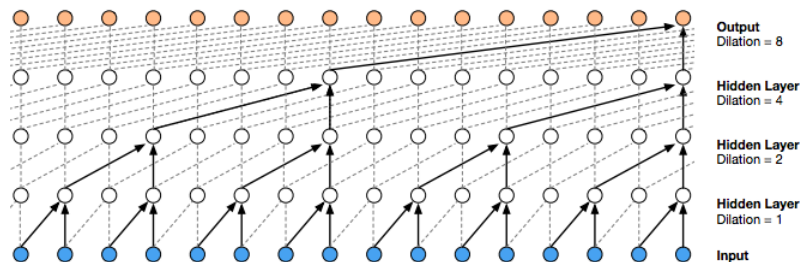
AUTOREGRESSIVE MODELS

TRAINING OBJECTIVE

Maximum likelihood = sum of conditional cross-entropies.

$$\max_{\theta} \mathbb{E}_{x \sim P} [\log p_{\theta}(x)] = \max_{\theta} \mathbb{E}_{x \sim P} \left[\sum_{i=1}^d \log p_{\theta}(x_i | x_{<i}) \right].$$

- ▶ Empirical objective: $-\frac{1}{N} \sum_n \sum_i \log p_{\theta}(x_{n,i} | x_{n,<i})$.
- ▶ *Teacher forcing*: condition on true prefixes $x_{<i}$ during training.
- ▶ Implementation: causal masking (sequences), masked convolutions (images), parallel loss over all positions.



AUTOREGRESSIVE MODELS

GENERATION (ANCESTRAL SAMPLING)

Ancestral sampling: draw variables one-by-one following the factorization order.

Algorithm 6: Autoregressive ancestral sampling

Input: learned conditionals $p_{\theta}(x_i \mid x_{<i})$, dimension d

Output: sample $x = (x_1, \dots, x_d)$

for $i \leftarrow 1$ **to** d **do**

 | Sample $x_i \sim p_{\theta}(\cdot \mid x_{<i})$

end

Notes:

- ▶ Exact and simple; sequential cost $\mathcal{O}(d)$ (limited parallelism at inference).
- ▶ For discrete outputs (text, pixels): categorical sampling; for continuous: Gaussian or mixture.
- ▶ Temperature/top- k /nucleus (p) sampling often used for text (heuristics, not MLE).

AUTOREGRESSIVE MODELS

LLMs AS AUTOREGRESSIVE NEXT-TOKEN PREDICTORS

Token sequence factorization. For tokens $w_{1:T}$,

$$p_{\theta}(w_{1:T}) = \prod_{t=1}^T p_{\theta}(w_t \mid w_{<t}), \quad \log p_{\theta}(w_{1:T}) = \sum_{t=1}^T \log p_{\theta}(w_t \mid w_{<t}).$$

- ▶ Transformer *decoder*-only with causal mask models $p_{\theta}(w_t \mid w_{<t})$.
- ▶ Training: minimize cross-entropy to true next token (teacher forcing, parallel over positions).
- ▶ Tokenization: subword units (BPE/WordPiece) turn text into discrete tokens; softmax over vocab.

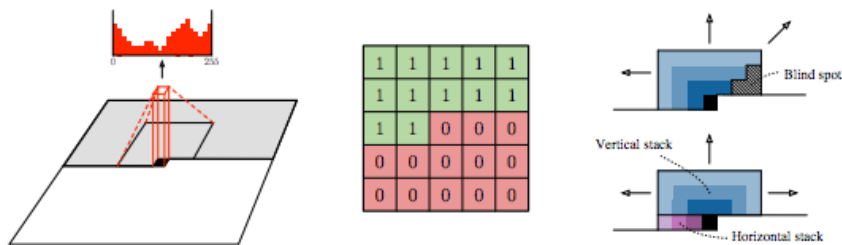
AUTOREGRESSIVE MODELS

PIXELCNN / PIXELRNN

Image factorization. Raster-scan ordering over pixels (and optionally channels):

$$p_{\theta}(x) = \prod_{i=1}^{HW \cdot C} p_{\theta}(x_i | x_{<i}).$$

- ▶ **Masked convolutions** enforce causality (no access to future pixels).
- ▶ **Receptive field** grows with layers; PixelRNN uses recurrent structure.
- ▶ Discrete pixels: categorical over 256 bins or mixture of logistics (PixelCNN++).



NORMALIZING FLOWS

OVERVIEW

A Normalizing Flow is usually seen as:

- ▶ a *generative model*,
- ▶ a *bijective mapping*,
- ▶ an *invertible neural network*,
- ▶ a *density estimator*.

NORMALIZING FLOWS

MAPPING BETWEEN DISTRIBUTIONS — POINT TO POINT

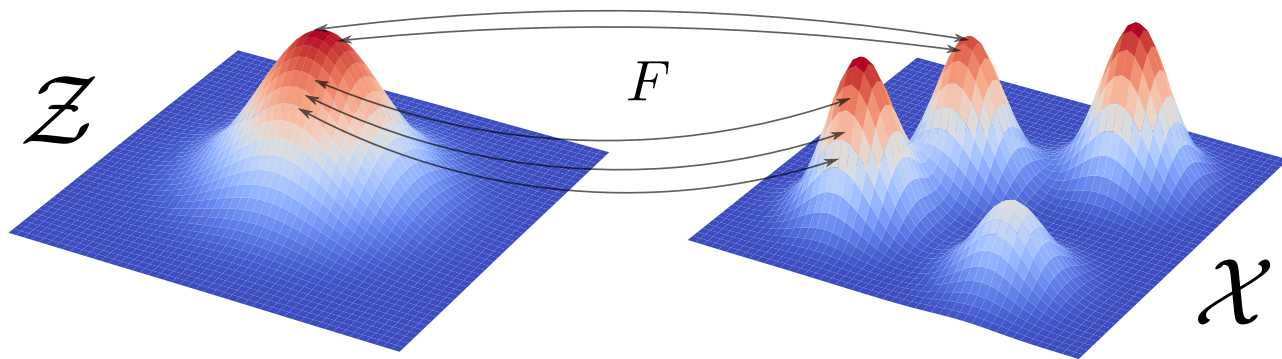


Figure. A mapping between two probability distributions
Point to point

NORMALIZING FLOWS

MAPPING BETWEEN DISTRIBUTIONS — SUBSET TO SUBSET

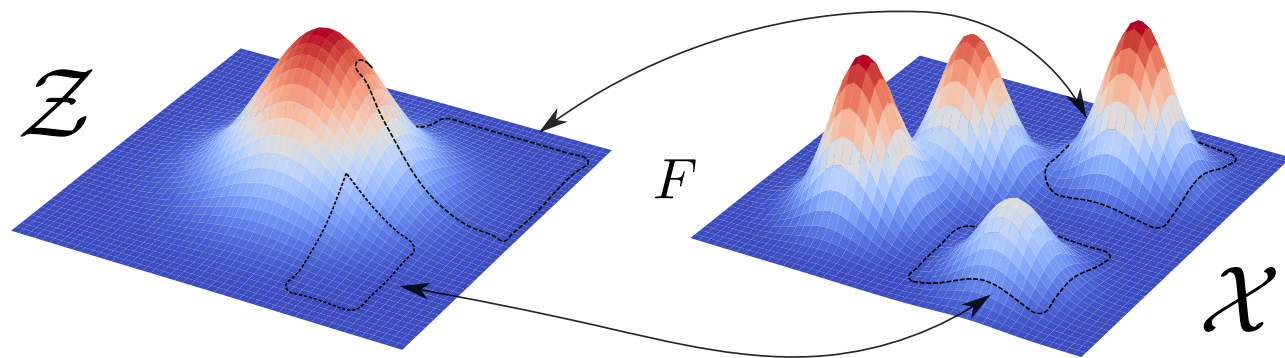


Figure. A mapping between two probability distributions
Subset to subset

NORMALIZING FLOWS

MATHEMATICAL FRAMEWORK

Normalizing Flow

A Normalizing Flow is a bijective function between a data space \mathcal{X} and a latent space \mathcal{Z} , both subsets of \mathbb{R}^d .

$$\begin{aligned} F : \mathcal{X} &\longmapsto \mathcal{Z} \\ x &\longmapsto z = F(x) \end{aligned}$$

Data and Latent Distributions

In theory, a NF maps a target distribution P (the data distribution) to a simple latent distribution Q . Usually, Q is set to be a multivariate normal $\mathcal{N}(0, \mathbf{I}_d)$. p and q denote the densities of P and Q respectively.

NORMALIZING FLOWS

HOW DOES IT WORK?

In practice, the mapping is *not perfect*. P induces a distribution Q through F , and the latent distribution Q induces \hat{P} through F^{-1} , which is the learned distribution. The forward pass F is called the *normalizing* direction while the inverse pass F^{-1} is called the *generative* direction.

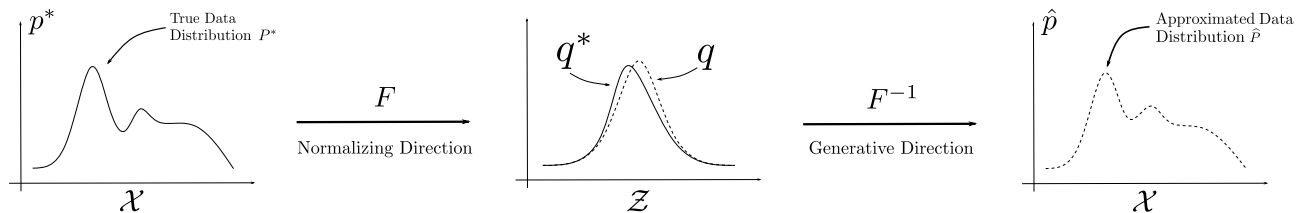


Figure. 1D Normalizing Flow process.

NORMALIZING FLOWS

CHANGE OF VARIABLES

Change of Variables Formula

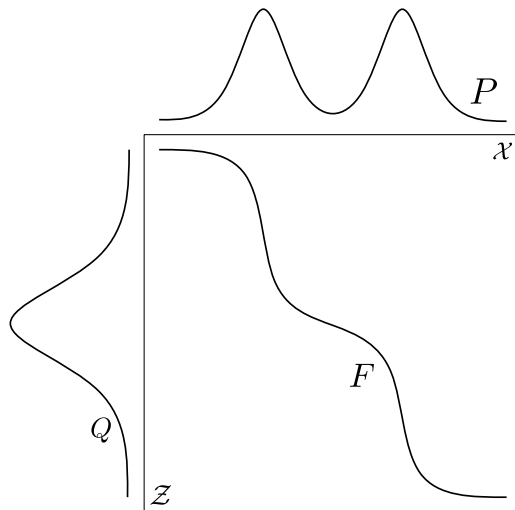
For a bijective and continuous function F and a latent distribution Q , the distribution induced by Q and F is defined by:

$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)). \quad (1)$$

NORMALIZING FLOWS

CHANGE OF VARIABLES — VISUAL INTUITION

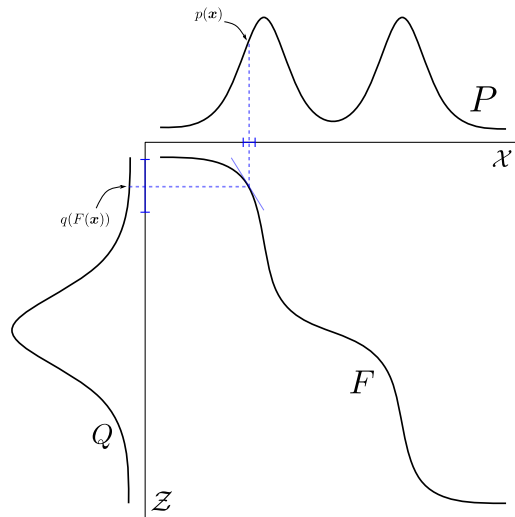
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



NORMALIZING FLOWS

CHANGE OF VARIABLES — VISUAL INTUITION

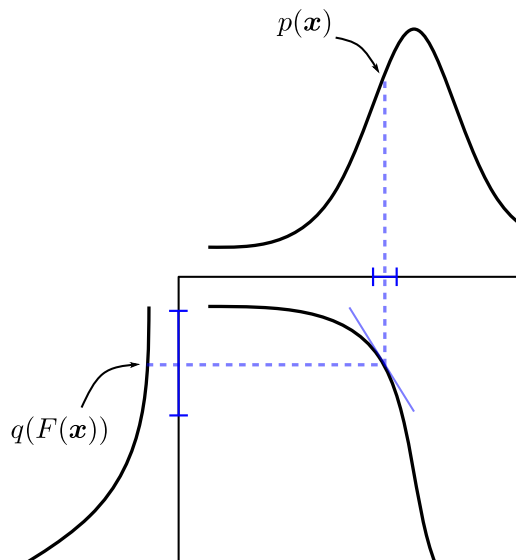
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



NORMALIZING FLOWS

CHANGE OF VARIABLES — VISUAL INTUITION

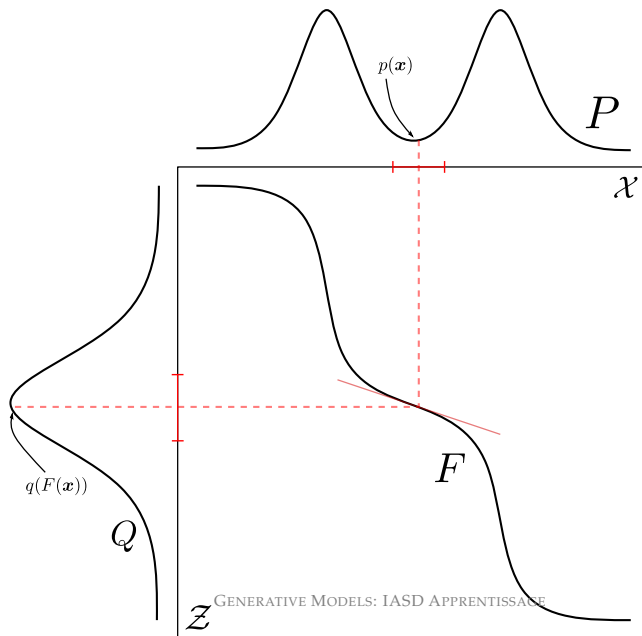
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



NORMALIZING FLOWS

CHANGE OF VARIABLES — VISUAL INTUITION

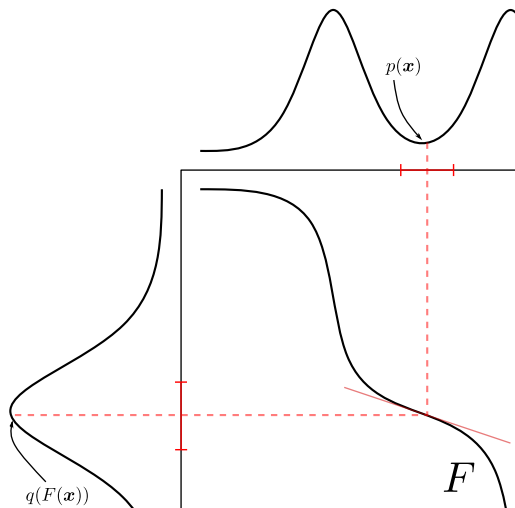
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



NORMALIZING FLOWS

CHANGE OF VARIABLES — VISUAL INTUITION

$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



NORMALIZING FLOWS

DENSITY ESTIMATION

To perform density estimation:

1. Draw $x \sim P$,
2. Compute $F(x)$ and $|\det \text{Jac}_F(x)|$,
3. Compute $\hat{p}(x) = q(F(x)) |\det \text{Jac}_F(x)|$.

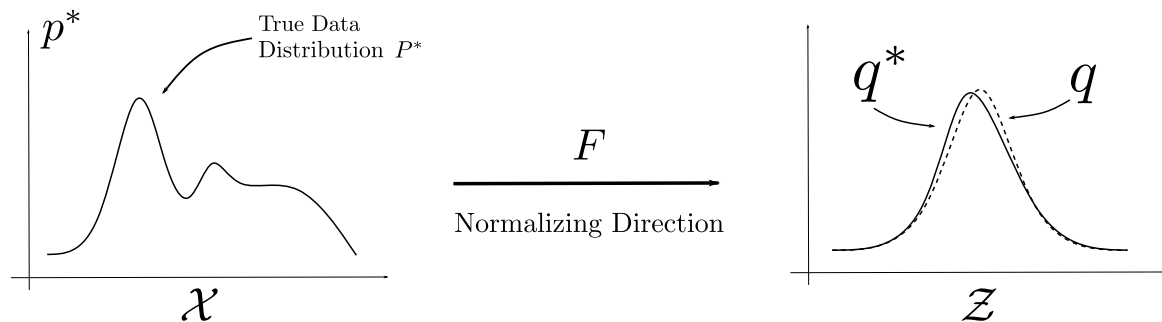


Figure. 1D Normalizing Flow process of density estimation.

NORMALIZING FLOWS

DATA GENERATION

To perform data generation:

1. Draw $z \sim Q$,
2. Compute $x = F^{-1}(z)$.

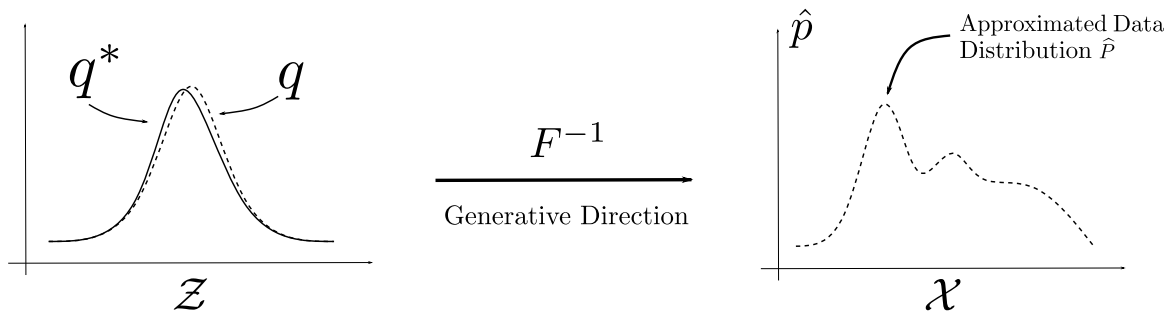


Figure. 1D Normalizing Flow process of generation.

NORMALIZING FLOWS

LEARNING STEPS

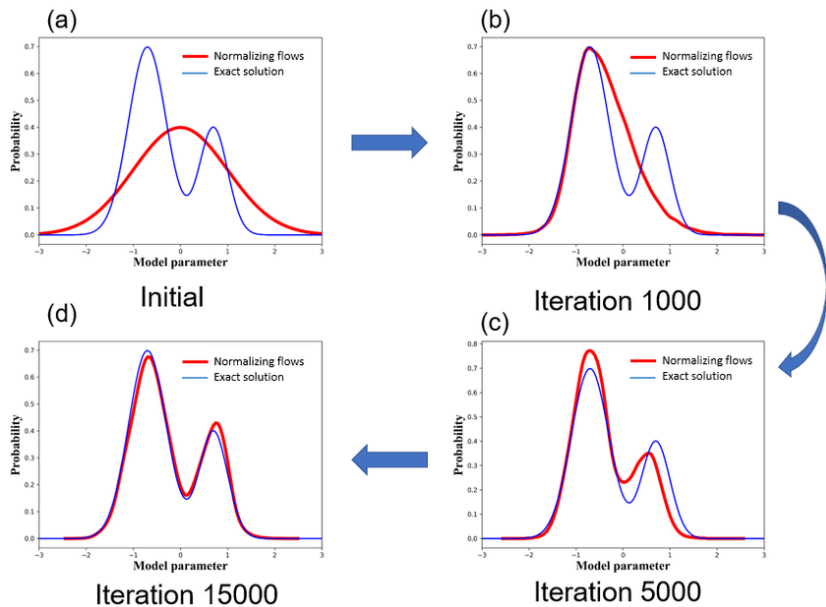


Figure. Learning process for a 1D Normalizing Flow.

ENERGY-BASED MODELS

PARTITION FUNCTION CHALLENGE

Definition. An Energy-Based Model (EBM) defines

$$p_{\theta}(x) = \frac{\exp(-E_{\theta}(x))}{Z(\theta)}, \quad Z(\theta) = \int \exp(-E_{\theta}(x)) dx.$$

- ▶ $E_{\theta}(x)$ is an *energy function* (low energy = high probability).
- ▶ $Z(\theta)$ is the *partition function* ensuring normalization.
- ▶ Problem: computing $Z(\theta)$ is generally intractable (high-dimensional integral).

Takeaway: normalization constant is the main pain in EBMs.

ENERGY-BASED MODELS

TRAINING — LOG-LIKELIHOOD AND GRADIENT

Log-likelihood for one sample x :

$$\log p_\theta(x) = -E_\theta(x) - \log Z(\theta), \quad \text{with } Z(\theta) = \int e^{-E_\theta(u)} du.$$

Gradient derivation:

$$\nabla_\theta \log p_\theta(x) = -\nabla_\theta E_\theta(x) - \nabla_\theta \log Z(\theta)$$

Compute $\nabla_\theta \log Z(\theta)$ explicitly:

$$\begin{aligned} Z(\theta) &= \int e^{-E_\theta(u)} du, & \nabla_\theta Z(\theta) &= \int e^{-E_\theta(u)} (-\nabla_\theta E_\theta(u)) du, \\ \nabla_\theta \log Z(\theta) &= \frac{1}{Z(\theta)} \nabla_\theta Z(\theta) \\ &= - \int \frac{e^{-E_\theta(u)}}{Z(\theta)} \nabla_\theta E_\theta(u) du \\ &= - \mathbb{E}_{u \sim p_\theta} [\nabla_\theta E_\theta(u)]. \end{aligned}$$

Substitute back:

$$\nabla_\theta \log p_\theta(x) = -\nabla_\theta E_\theta(x) - \nabla_\theta \log Z(\theta) = -\nabla_\theta E_\theta(x) + \mathbb{E}_{u \sim p_\theta} [\nabla_\theta E_\theta(u)].$$

ENERGY-BASED MODELS

TRAINING — PRACTICAL GRADIENT STEP

Empirical objective (dataset $\{x_n\}$):

$$\nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(x_n) = -\frac{1}{N} \sum_{n=1}^N \nabla_{\theta} E_{\theta}(x_n) + \mathbb{E}_{u \sim p_{\theta}} [\nabla_{\theta} E_{\theta}(u)].$$

Gradient step (schematic):

$$\theta \leftarrow \theta - \eta \left(-\frac{1}{N} \sum_n \nabla_{\theta} E_{\theta}(x_n) + \mathbb{E}_{u \sim p_{\theta}} [\nabla_{\theta} E_{\theta}(u)] \right).$$

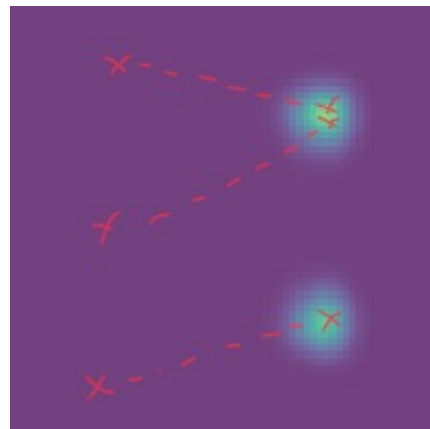
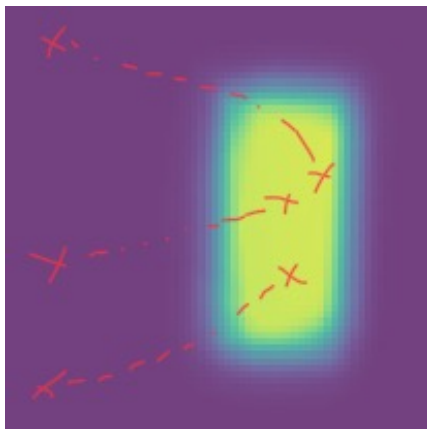
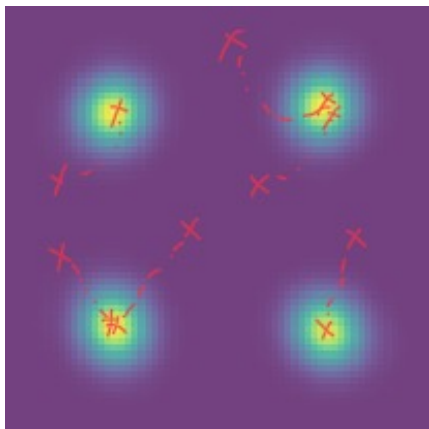
- ▶ *Pull down* energy on data (first term), *push up* on model samples (second term).
- ▶ Trade-offs: bias vs. mixing time; stability tricks (noise scale, step size, gradient clipping, spectral norm).

ENERGY-BASED MODELS

SAMPLING

Goal: generate $x \sim p_\theta(x) \propto e^{-E_\theta(x)}$.

- ▶ Direct sampling is impossible (requires $Z(\theta)$).
- ▶ Use MCMC methods (e.g., Langevin dynamics, Hamiltonian Monte Carlo).
- ▶ Iteratively update $x \leftarrow x - \eta \nabla_x E_\theta(x) + \sqrt{2\eta} \xi$, $\xi \sim \mathcal{N}(0, I)$.
- ▶ Paths follow the energy landscape toward low-energy regions (data modes).



AUTOENCODER

BASIC ARCHITECTURE AND OBJECTIVE

Architecture:

- ▶ **Encoder:** maps input x to a low-dimensional latent representation $z = f_\phi(x)$.
- ▶ **Decoder:** reconstructs input from latent z , i.e., $\hat{x} = g_\theta(z)$.
- ▶ The model is trained end-to-end to minimize the difference between x and \hat{x} .

Objective:

$$\min_{\theta, \phi} \mathbb{E}_{x \sim P_{\text{data}}} [\ell(x, g_\theta(f_\phi(x)))]$$

where ℓ is typically mean squared error: $\ell(x, \hat{x}) = \|x - \hat{x}\|^2$. **Limitations for generative modeling:**

- ▶ No explicit generative process for sampling new data from the latent space.
- ▶ Latent space may not follow a known distribution—sampling z at random often yields unrealistic outputs.
- ▶ Not a true probabilistic model; lacks explicit likelihood or regularization of latent space.

AUTOENCODER

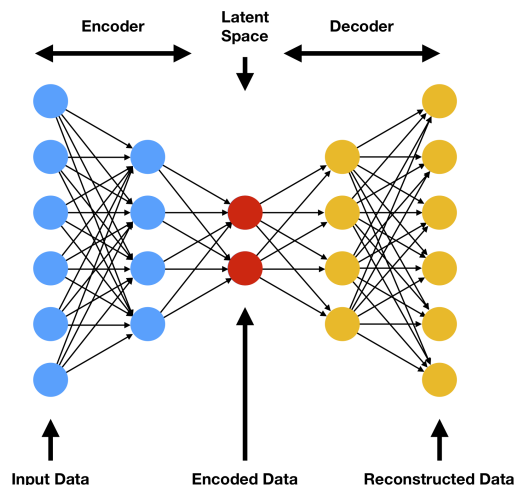
ILLUSTRATION AND EXAMPLE

Architecture:

- ▶ **Input:** x (e.g., image, signal)
- ▶ **Encoder:** compresses x to latent z
- ▶ **Decoder:** reconstructs \hat{x} from z

Applications:

- ▶ Dimensionality reduction, denoising, anomaly detection, feature learning
- ▶ Not directly suited for generating novel samples



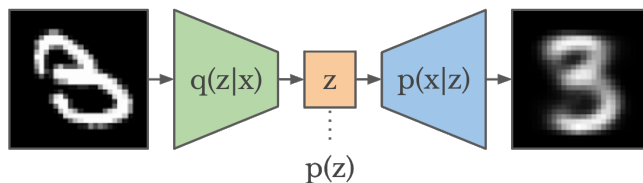
VARIATIONAL AUTOENCODER

INTUITION BEHIND VAE

Key ideas:

- ▶ **Probabilistic encoder:** Instead of mapping $x \rightarrow z$ deterministically, encode x as a *distribution* over latent variables: $q_\phi(z|x)$ (e.g., Gaussian with mean and variance predicted by encoder).
- ▶ **Probabilistic decoder:** Model $p_\theta(x|z)$, i.e., generate x from latent z .
- ▶ **Latent variable modeling:** Place a prior $p(z)$ (usually standard normal) on the latent space to encourage structure and enable sampling.
- ▶ **Regularization:** Use KL divergence $\mathcal{D}_{\text{KL}}(q_\phi(z|x) \| p(z))$ to encourage $q_\phi(z|x)$ to be close to the prior, making the latent space well-behaved and suitable for generative sampling.

Summary: VAE is a probabilistic autoencoder that learns both to reconstruct data and to regularize the latent space for generative use.



VARIATIONAL AUTOENCODER

EVIDENCE LOWER BOUND (ELBO)

Objective: Evidence Lower Bound (ELBO)

$$\log p_{\theta}(x) \geq \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - \mathcal{D}_{\text{KL}}(q_{\phi}(z|x) || p(z))$$

- ▶ **First term:** Expected log-likelihood (reconstruction accuracy)
- ▶ **Second term:** KL divergence between encoder distribution and prior (regularization)

Relation to MLE: Maximizing ELBO approximates maximizing the marginal likelihood $p_{\theta}(x)$ (i.e., maximum likelihood estimation for latent variable models).

Training strategy:

- ▶ Optimize the ELBO jointly with respect to encoder (ϕ) and decoder (θ) parameters.
- ▶ Use stochastic gradient descent with the reparameterization trick to backpropagate through stochastic nodes.

VARIATIONAL AUTOENCODER

REPARAMETERIZATION TRICK

Intuition:

- ▶ Allows gradients to flow through random sampling by expressing sampling as a deterministic function of parameters and noise.
- ▶ Enables efficient and low-variance gradient estimation for stochastic variables.

Mathematical formulation:

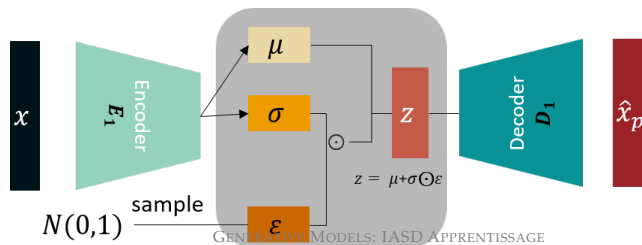
- ▶ For $q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \sigma_\phi(x)^2)$,
- ▶ Sample z as:

$$z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

- ▶ Now, z is a deterministic function of x , ϕ , and random noise ϵ .

Benefits:

- ▶ Enables backpropagation through stochastic sampling.
- ▶ Crucial for training VAEs with gradient-based methods.



VARIATIONAL AUTOENCODER

LOSS

VAE Loss Function:

$$\mathcal{L}_{\text{VAE}}(x) = \underbrace{\mathbb{E}_{q_{\phi}(z|x)}[-\log p_{\theta}(x|z)]}_{\text{Reconstruction loss}} + \underbrace{\mathcal{D}_{\text{KL}}(q_{\phi}(z|x) || p(z))}_{\text{KL divergence (regularization)}}$$

- ▶ **Reconstruction loss:** Measures how well the decoder can reconstruct the input from the latent code.
- ▶ **KL divergence:** Encourages the approximate posterior $q_{\phi}(z|x)$ to match the prior $p(z)$ (e.g., standard normal), regularizing the latent space.
- ▶ **Trade-off:** Balances data fidelity (reconstruction) and latent space regularity (generative quality). Too much weight on KL: blurry reconstructions; too little: latent space collapse.

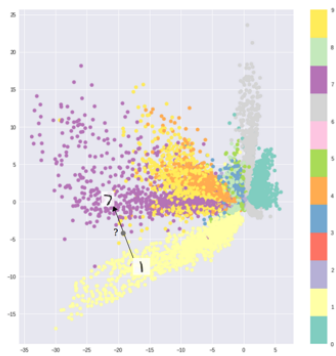
Typical formula:

$$\mathcal{L}_{\text{VAE}}(x) = \frac{1}{2} \sum_j \left(\sigma_j^2(x) + \mu_j(x)^2 - 1 - \log \sigma_j^2(x) \right) + \text{Reconstruction loss}$$

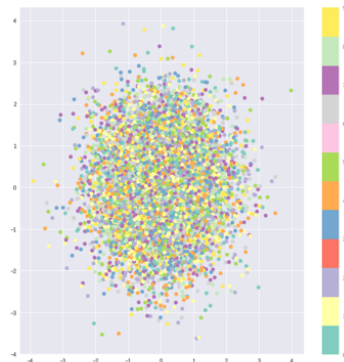
VARIATIONAL AUTOENCODER

Loss

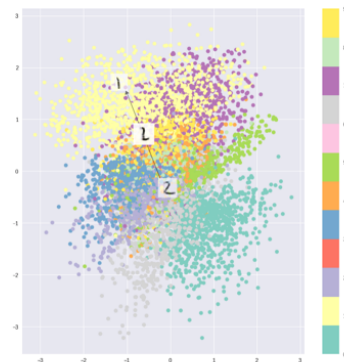
Only reconstruction loss



Only KL divergence



Combination



VARIANTS OF VAE

EXTENSIONS AND IMPROVEMENTS

Key VAE Variants:

- ▶ **Conditional VAE (CVAE):** Conditions both encoder and decoder on auxiliary information (e.g., labels, attributes) to enable conditional generation [5].
- ▶ **β -VAE:** Introduces a hyperparameter β to scale the KL term, encouraging disentangled latent representations [2].
- ▶ **Other notable extensions:**
 - *VampPrior*: Learnable mixture prior for more flexible latent space.
 - *Vector Quantized VAE (VQ-VAE)* [4]: Discrete latent variables via vector quantization.
 - *Hierarchical VAE*: Multiple layers of latent variables.
 - *FactorVAE*, *InfoVAE*, *WAE* (Wasserstein Autoencoder), etc.

References:

- ▶ Original VAE paper: [3]
- ▶ Conditional VAE: [5]
- ▶ β -VAE : [2]
- ▶ VQ-VAE: [4]

VQ-VAE (VAN DEN OORD ET AL., 2017)

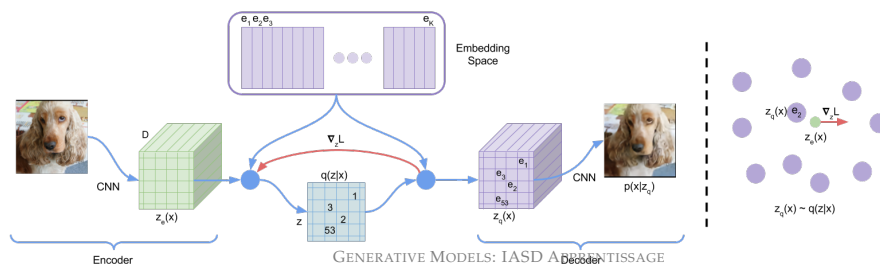
VECTOR QUANTIZED VAE

Key ideas:

- ▶ Introduces **discrete latent representations** via *vector quantization*—the encoder outputs are mapped to the nearest entry in a learned codebook.
- ▶ The decoder reconstructs x from the quantized latent code.
- ▶ Enables modeling of discrete structure in data (e.g., language, audio, images).
- ▶ Discrete latents are particularly beneficial for combining VAEs with powerful generative models (such as PixelCNN or Transformers) in the latent space.
- ▶ Facilitates improved sample quality and more interpretable representations.

Benefits:

- ▶ Enables use of GAN-like or autoregressive models in discrete latent space.
- ▶ Improved performance on high-fidelity image and audio generation tasks.



REFERENCES I

- [1] Kate Bibbings, Peter J. Harding, Ian D. Loram, Nicholas Combes, and Emma F. Hodson-Tole. Foreground Detection Analysis of Ultrasound Image Sequences Identifies Markers of Motor Neurone Disease across Diagnostically Relevant Skeletal Muscles. *Ultrasound in Medicine & Biology*, 45(5):1164–1175, May 2019. ISSN 1879-291X. doi: 10.1016/j.ultrasmedbio.2019.01.018.
- [2] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: LEARNING BASIC VISUAL CONCEPTS WITH A CONSTRAINED VARIATIONAL FRAMEWORK. 2017.
- [3] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes, December 2022. URL <http://arxiv.org/abs/1312.6114>. arXiv:1312.6114 [cs, stat].
- [4] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation Learning, May 2018. URL <http://arxiv.org/abs/1711.00937>. arXiv:1711.00937 [cs].
- [5] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning Structured Output Representation using Deep Conditional Generative Models. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://papers.nips.cc/paper_files/paper/2015/hash/8d55a249e6baa5c06772297520da2051-Abstract.html.