

Deep Generative Models

Fall 2025

Sharif University of Technology



Alireza Mirrokni - 401106617

Problem Set 1

Problem 1: Properties of Gaussian Distributions

We first find the Moment Generating Function (MGF) of a Gaussian random vector. Let $\mathbf{X} \in \mathbb{R}^d$ be Gaussian with mean $\boldsymbol{\mu}$ and covariance Σ . Choose a matrix $L \in \mathbb{R}^{d \times d}$ such that $LL^\top = \Sigma$ and let $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, I_d)$ with independent coordinates. Then \mathbf{X} has the representation $\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + L\mathbf{Z}$. For any $\mathbf{t} \in \mathbb{R}^d$,

$$\begin{aligned} M_{\mathbf{X}}(\mathbf{t}) &= \mathbb{E} [\exp(\mathbf{t}^\top \mathbf{X})] = \mathbb{E} [\exp(\mathbf{t}^\top (\boldsymbol{\mu} + L\mathbf{Z}))] \\ &= \exp(\mathbf{t}^\top \boldsymbol{\mu}) \mathbb{E} [\exp((L^\top \mathbf{t})^\top \mathbf{Z})]. \end{aligned}$$

Let $\mathbf{u} = L^\top \mathbf{t}$. Since Z_1, \dots, Z_d are independent standard normals,

$$\mathbb{E} [\exp(\mathbf{u}^\top \mathbf{Z})] = \mathbb{E} \left[\exp \left(\sum_{i=1}^d u_i Z_i \right) \right] = \prod_{i=1}^d \mathbb{E} [\exp(u_i Z_i)].$$

For scalar $Z \sim \mathcal{N}(0, 1)$ and $u \in \mathbb{R}$,

$$\begin{aligned} \mathbb{E} [\exp(uZ)] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp \left(uz - \frac{1}{2} z^2 \right) dz \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp \left(-\frac{1}{2} (z-u)^2 + \frac{1}{2} u^2 \right) dz = \exp \left(\frac{1}{2} u^2 \right), \end{aligned}$$

because $\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp \left(-\frac{1}{2} (z-u)^2 \right) dz = 1$. Therefore

$$\mathbb{E} [\exp(\mathbf{u}^\top \mathbf{Z})] = \prod_{i=1}^d \exp \left(\frac{1}{2} u_i^2 \right) = \exp \left(\frac{1}{2} \mathbf{u}^\top \mathbf{u} \right).$$

Substituting $\mathbf{u} = L^\top \mathbf{t}$ yields

$$M_{\mathbf{X}}(\mathbf{t}) = \exp(\mathbf{t}^\top \boldsymbol{\mu}) \exp \left(\frac{1}{2} \mathbf{t}^\top LL^\top \mathbf{t} \right) = \exp \left(\mathbf{t}^\top \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^\top \Sigma \mathbf{t} \right).$$

1. For any $\mathbf{s} \in \mathbb{R}^m$,

$$\begin{aligned}
M_{\mathbf{Y}}(\mathbf{s}) &= \mathbb{E} \left[\exp \left(\mathbf{s}^\top \mathbf{Y} \right) \right] = \mathbb{E} \left[\exp \left(\mathbf{s}^\top (A\mathbf{X} + \mathbf{b}) \right) \right] \\
&= \exp \left(\mathbf{s}^\top \mathbf{b} \right) \mathbb{E} \left[\exp \left((A^\top \mathbf{s})^\top \mathbf{X} \right) \right] = \exp \left(\mathbf{s}^\top \mathbf{b} \right) M_{\mathbf{X}}(A^\top \mathbf{s}) \\
&= \exp \left(\mathbf{s}^\top \mathbf{b} + (A^\top \mathbf{s})^\top \boldsymbol{\mu} + \frac{1}{2} (A^\top \mathbf{s})^\top \boldsymbol{\Sigma} (A^\top \mathbf{s}) \right) \\
&= \exp \left(\mathbf{s}^\top (A\boldsymbol{\mu} + \mathbf{b}) + \frac{1}{2} \mathbf{s}^\top (A\boldsymbol{\Sigma} A^\top) \mathbf{s} \right).
\end{aligned}$$

This is the mgf of a multivariate normal distribution with mean $A\boldsymbol{\mu} + \mathbf{b}$ and covariance $A\boldsymbol{\Sigma} A^\top$. By uniqueness of the mgf, \mathbf{Y} is Gaussian with

$$\mathbb{E}[\mathbf{Y}] = A\boldsymbol{\mu} + \mathbf{b}, \quad \text{Cov}(\mathbf{Y}) = A\boldsymbol{\Sigma} A^\top.$$

Equivalently, the cumulant generating function $K_{\mathbf{Y}}(\mathbf{s}) = \log M_{\mathbf{Y}}(\mathbf{s}) = \mathbf{s}^\top (A\boldsymbol{\mu} + \mathbf{b}) + \frac{1}{2} \mathbf{s}^\top (A\boldsymbol{\Sigma} A^\top) \mathbf{s}$ satisfies

$$\nabla_s K_{\mathbf{Y}}(\mathbf{0}) = A\boldsymbol{\mu} + \mathbf{b}, \quad \nabla_s^2 K_{\mathbf{Y}}(\mathbf{0}) = A\boldsymbol{\Sigma} A^\top.$$

2. For any $\mathbf{s} \in \mathbb{R}^d$,

$$\begin{aligned}
M_{\mathbf{Z}}(\mathbf{s}) &= \mathbb{E} \left[\exp \left(\mathbf{s}^\top \mathbf{Z} \right) \right] = \mathbb{E} \left[\exp \left(\mathbf{s}^\top (a\mathbf{X} + b\mathbf{Y}) \right) \right] \\
&= \mathbb{E} \left[\exp \left(a\mathbf{s}^\top \mathbf{X} \right) \exp \left(b\mathbf{s}^\top \mathbf{Y} \right) \right] = \mathbb{E} \left[\exp \left(a\mathbf{s}^\top \mathbf{X} \right) \right] \mathbb{E} \left[\exp \left(b\mathbf{s}^\top \mathbf{Y} \right) \right] \\
&= M_{\mathbf{X}}(a\mathbf{s}) M_{\mathbf{Y}}(b\mathbf{s}) \\
&= \exp \left((a\mathbf{s})^\top \boldsymbol{\mu}_X + \frac{1}{2} (a\mathbf{s})^\top \boldsymbol{\Sigma}_X (a\mathbf{s}) \right) \exp \left((b\mathbf{s})^\top \boldsymbol{\mu}_Y + \frac{1}{2} (b\mathbf{s})^\top \boldsymbol{\Sigma}_Y (b\mathbf{s}) \right) \\
&= \exp \left(\mathbf{s}^\top (a\boldsymbol{\mu}_X + b\boldsymbol{\mu}_Y) + \frac{1}{2} \mathbf{s}^\top (a^2 \boldsymbol{\Sigma}_X + b^2 \boldsymbol{\Sigma}_Y) \mathbf{s} \right).
\end{aligned}$$

This is the mgf of a multivariate normal distribution with mean $a\boldsymbol{\mu}_X + b\boldsymbol{\mu}_Y$ and covariance $a^2 \boldsymbol{\Sigma}_X + b^2 \boldsymbol{\Sigma}_Y$. Again by uniqueness of the mgf, \mathbf{Z} is Gaussian with

$$\boldsymbol{\mu}_{\mathbf{Z}} = a\boldsymbol{\mu}_X + b\boldsymbol{\mu}_Y, \quad \boldsymbol{\Sigma}_{\mathbf{Z}} = a^2 \boldsymbol{\Sigma}_X + b^2 \boldsymbol{\Sigma}_Y.$$

Equivalently, the cumulant generating function $K_{\mathbf{Z}}(\mathbf{s}) = \log M_{\mathbf{Z}}(\mathbf{s}) = \mathbf{s}^\top (a\boldsymbol{\mu}_X + b\boldsymbol{\mu}_Y) + \frac{1}{2} \mathbf{s}^\top (a^2 \boldsymbol{\Sigma}_X + b^2 \boldsymbol{\Sigma}_Y) \mathbf{s}$ satisfies

$$\nabla_s K_{\mathbf{Z}}(\mathbf{0}) = a\boldsymbol{\mu}_X + b\boldsymbol{\mu}_Y, \quad \nabla_s^2 K_{\mathbf{Z}}(\mathbf{0}) = a^2 \boldsymbol{\Sigma}_X + b^2 \boldsymbol{\Sigma}_Y.$$

3. For any $\mathbf{s} \in \mathbb{R}^{d_y}$,

$$M_{\mathbf{Y}}(\mathbf{s}) = \mathbb{E} \left[\exp \left(\mathbf{s}^\top \mathbf{Y} \right) \right] = \mathbb{E} \left[\mathbb{E} \left[\exp \left(\mathbf{s}^\top \mathbf{Y} \right) \mid \mathbf{X} \right] \right].$$

Given \mathbf{X} , the vector \mathbf{Y} is Gaussian with mean $A\mathbf{X} + \mathbf{b}$ and covariance $\boldsymbol{\Sigma}_{Y|X}$, hence its conditional mgf is

$$\mathbb{E} \left[\exp \left(\mathbf{s}^\top \mathbf{Y} \right) \mid \mathbf{X} \right] = \exp \left(\mathbf{s}^\top (A\mathbf{X} + \mathbf{b}) + \frac{1}{2} \mathbf{s}^\top \boldsymbol{\Sigma}_{Y|X} \mathbf{s} \right).$$

Taking expectation over \mathbf{X} and using the Gaussian mgf for \mathbf{X} at argument $A^\top \mathbf{s}$ gives

$$\begin{aligned} M_{\mathbf{Y}}(\mathbf{s}) &= \exp\left(\mathbf{s}^\top \mathbf{b} + \frac{1}{2}\mathbf{s}^\top \Sigma_{Y|X}\mathbf{s}\right) \mathbb{E}\left[\exp\left(\left(A^\top \mathbf{s}\right)^\top \mathbf{X}\right)\right] \\ &= \exp\left(\mathbf{s}^\top \mathbf{b} + \frac{1}{2}\mathbf{s}^\top \Sigma_{Y|X}\mathbf{s}\right) \exp\left(\left(A^\top \mathbf{s}\right)^\top \boldsymbol{\mu}_X + \frac{1}{2}\left(A^\top \mathbf{s}\right)^\top \Sigma_{XX}\left(A^\top \mathbf{s}\right)\right) \\ &= \exp\left(\mathbf{s}^\top (A\boldsymbol{\mu}_X + \mathbf{b}) + \frac{1}{2}\mathbf{s}^\top \left(\Sigma_{Y|X} + A\Sigma_{XX}A^\top\right)\mathbf{s}\right). \end{aligned}$$

This is the mgf of a multivariate normal distribution with mean $\boldsymbol{\mu}_Y = A\boldsymbol{\mu}_X + \mathbf{b}$ and covariance $\Sigma_Y = \Sigma_{Y|X} + A\Sigma_{XX}A^\top$. By uniqueness of the mgf, the marginal \mathbf{Y} is Gaussian with

$$\boldsymbol{\mu}_Y = A\boldsymbol{\mu}_X + \mathbf{b}, \quad \Sigma_Y = \Sigma_{Y|X} + A\Sigma_{XX}A^\top.$$

As a consistency check, the same formulas follow from the laws of total expectation and total covariance:

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathbb{E}[\mathbf{Y} | \mathbf{X}]] = \mathbb{E}[A\mathbf{X} + \mathbf{b}] = A\boldsymbol{\mu}_X + \mathbf{b},$$

$$\begin{aligned} \text{Cov}(\mathbf{Y}) &= \mathbb{E}[\text{Cov}(\mathbf{Y} | \mathbf{X})] + \text{Cov}(\mathbb{E}[\mathbf{Y} | \mathbf{X}]) \\ &= \Sigma_{Y|X} + \text{Cov}(A\mathbf{X} + \mathbf{b}) = \Sigma_{Y|X} + A\Sigma_{XX}A^\top. \end{aligned}$$

4. Write $\delta\mathbf{x} = \mathbf{x} - \boldsymbol{\mu}_X$ and $\delta\mathbf{y} = \mathbf{y} - \boldsymbol{\mu}_Y$. The joint density is

$$f_{\mathbf{X},\mathbf{Y}}(\mathbf{x},\mathbf{y}) = (2\pi)^{-\frac{d_x+d_y}{2}} (\det \Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\begin{pmatrix} \delta\mathbf{x} \\ \delta\mathbf{y} \end{pmatrix}^\top \Sigma^{-1} \begin{pmatrix} \delta\mathbf{x} \\ \delta\mathbf{y} \end{pmatrix}\right).$$

The marginal density of \mathbf{X} is

$$f_{\mathbf{X}}(\mathbf{x}) = (2\pi)^{-\frac{d_x}{2}} (\det \Sigma_{XX})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\delta\mathbf{x}^\top \Sigma_{XX}^{-1} \delta\mathbf{x}\right).$$

Hence

$$f_{\mathbf{Y}|\mathbf{X}=\mathbf{x}}(\mathbf{y}) = \frac{f_{\mathbf{X},\mathbf{Y}}(\mathbf{x},\mathbf{y})}{f_{\mathbf{X}}(\mathbf{x})} = (2\pi)^{-\frac{d_y}{2}} \left(\frac{\det \Sigma_{XX}}{\det \Sigma}\right)^{\frac{1}{2}} \exp\left(-\frac{1}{2}Q(\mathbf{x},\mathbf{y})\right),$$

where

$$Q(\mathbf{x},\mathbf{y}) = \begin{pmatrix} \delta\mathbf{x} \\ \delta\mathbf{y} \end{pmatrix}^\top \Sigma^{-1} \begin{pmatrix} \delta\mathbf{x} \\ \delta\mathbf{y} \end{pmatrix} - \delta\mathbf{x}^\top \Sigma_{XX}^{-1} \delta\mathbf{x}.$$

Introduce the Schur complement

$$\Sigma_{Y|X} = \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}.$$

By the block matrix inversion identity proved below,

$$\Sigma^{-1} = \begin{pmatrix} \Sigma_{XX}^{-1} + \Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{Y|X}^{-1}\Sigma_{YX}\Sigma_{XX}^{-1} & -\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{Y|X}^{-1} \\ -\Sigma_{Y|X}^{-1}\Sigma_{YX}\Sigma_{XX}^{-1} & \Sigma_{Y|X}^{-1} \end{pmatrix}.$$

Substituting this into $Q(\mathbf{x}, \mathbf{y})$ and expanding gives

$$\begin{aligned} Q(\mathbf{x}, \mathbf{y}) &= \delta \mathbf{y}^\top \Sigma_{Y|X}^{-1} \delta \mathbf{y} - 2\delta \mathbf{y}^\top \Sigma_{Y|X}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} \delta \mathbf{x} + \delta \mathbf{x}^\top \Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{Y|X}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} \delta \mathbf{x} \\ &= (\delta \mathbf{y} - \Sigma_{YX} \Sigma_{XX}^{-1} \delta \mathbf{x})^\top \Sigma_{Y|X}^{-1} (\delta \mathbf{y} - \Sigma_{YX} \Sigma_{XX}^{-1} \delta \mathbf{x}). \end{aligned}$$

Therefore with $\mathbf{u} = \mathbf{y} - \boldsymbol{\mu}_Y - \Sigma_{YX} \Sigma_{XX}^{-1} (\mathbf{x} - \boldsymbol{\mu}_X)$,

$$f_{\mathbf{Y}|\mathbf{X}=\mathbf{x}}(\mathbf{y}) = (2\pi)^{-\frac{d_y}{2}} \left(\frac{\det \Sigma_{XX}}{\det \Sigma} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2} \mathbf{u}^\top \Sigma_{Y|X}^{-1} \mathbf{u} \right).$$

Using the block determinant identity $\det \Sigma = \det \Sigma_{XX} \det(\Sigma_{Y|X})$ proved below, the leading constant simplifies to $(2\pi)^{-\frac{d_y}{2}} (\det \Sigma_{Y|X})^{-\frac{1}{2}}$. Hence the conditional distribution is Gaussian with

$$\mathbf{Y} | \mathbf{X} = \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (\mathbf{x} - \boldsymbol{\mu}_X), \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}).$$

1. Proof of the block matrix inversion identity:

Assume Σ_{XX} is invertible and set $\Sigma_{Y|X} = \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$. Factor Σ via block Gaussian elimination

$$\Sigma = \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix} = \underbrace{\begin{pmatrix} I & 0 \\ \Sigma_{YX} \Sigma_{XX}^{-1} & I \end{pmatrix}}_L \underbrace{\begin{pmatrix} \Sigma_{XX} & 0 \\ 0 & \Sigma_{Y|X} \end{pmatrix}}_D \underbrace{\begin{pmatrix} I & \Sigma_{XX}^{-1} \Sigma_{XY} \\ 0 & I \end{pmatrix}}_U.$$

Inverting the triangular factors,

$$L^{-1} = \begin{pmatrix} I & 0 \\ -\Sigma_{YX} \Sigma_{XX}^{-1} & I \end{pmatrix}, \quad D^{-1} = \begin{pmatrix} \Sigma_{XX}^{-1} & 0 \\ 0 & \Sigma_{Y|X}^{-1} \end{pmatrix}, \quad U^{-1} = \begin{pmatrix} I & -\Sigma_{XX}^{-1} \Sigma_{XY} \\ 0 & I \end{pmatrix}.$$

Thus $\Sigma^{-1} = U^{-1} D^{-1} L^{-1}$. A direct multiplication gives

$$\Sigma^{-1} = \begin{pmatrix} \Sigma_{XX}^{-1} + \Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{Y|X}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} & -\Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{Y|X}^{-1} \\ -\Sigma_{Y|X}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} & \Sigma_{Y|X}^{-1} \end{pmatrix}.$$

2. Proof of the block determinant identity:

With the same L, D, U factorization as above, both L and U are block unit triangular, hence $\det L = 1$ and $\det U = 1$. Therefore

$$\begin{aligned} \det \Sigma &= \det L \det D \det U = \det D = \det \Sigma_{XX} \det(\Sigma_{Y|X}) \\ &= \det \Sigma_{XX} \det(\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}). \end{aligned}$$

5.

$$p_i(\mathbf{x}) = (2\pi)^{-d/2} (\det \Sigma_i)^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right), \quad i \in \{0, 1\}.$$

By definition,

$$D_{KL}(\mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \| \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)) = \mathbb{E}_{p_0} [\log p_0(\mathbf{X}) - \log p_1(\mathbf{X})].$$

Using the Gaussian log density and cancelling the $(-\frac{d}{2} \log(2\pi))$ terms,

$$\log p_0(\mathbf{x}) - \log p_1(\mathbf{x}) = \frac{1}{2} \log \frac{\det \Sigma_1}{\det \Sigma_0} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_0)^\top \Sigma_0^{-1} (\mathbf{x} - \boldsymbol{\mu}_0) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1).$$

Taking expectation under $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$,

$$D_{KL} = \frac{1}{2} \log \frac{\det \Sigma_1}{\det \Sigma_0} - \frac{1}{2} \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_0)^\top \Sigma_0^{-1} (\mathbf{X} - \boldsymbol{\mu}_0) \right] + \frac{1}{2} \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1} (\mathbf{X} - \boldsymbol{\mu}_1) \right].$$

The first expectation equals d because

$$\mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_0) (\mathbf{X} - \boldsymbol{\mu}_0)^\top \right] = \Sigma_0 \implies \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_0)^\top \Sigma_0^{-1} (\mathbf{X} - \boldsymbol{\mu}_0) \right] = \text{tr}(\Sigma_0^{-1} \Sigma_0) = d.$$

For the second expectation, write

$$\begin{aligned} \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_1) (\mathbf{X} - \boldsymbol{\mu}_1)^\top \right] &= \mathbb{E} \left[((\mathbf{X} - \boldsymbol{\mu}_0) + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)) ((\mathbf{X} - \boldsymbol{\mu}_0) + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1))^\top \right] \\ &= \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_0) (\mathbf{X} - \boldsymbol{\mu}_0)^\top \right] + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \\ &= \Sigma_0 + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top, \end{aligned}$$

since $\mathbb{E}[\mathbf{X} - \boldsymbol{\mu}_0] = \mathbf{0}$ kills the cross terms. Hence

$$\begin{aligned} \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1} (\mathbf{X} - \boldsymbol{\mu}_1) \right] &= \text{tr} \left(\Sigma_1^{-1} \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_1) (\mathbf{X} - \boldsymbol{\mu}_1)^\top \right] \right) \\ &= \text{tr}(\Sigma_1^{-1} \Sigma_0) + \text{tr} \left(\Sigma_1^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \right) \\ &= \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0), \end{aligned}$$

where the last step uses $\text{tr}(AB) = \text{tr}(BA)$ and $\text{tr}(\mathbf{a}\mathbf{a}^\top B) = \mathbf{a}^\top B\mathbf{a}$. Substituting,

$$\begin{aligned} D_{KL}(\mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \| \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)) &= \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - d + \log \frac{\det \Sigma_1}{\det \Sigma_0} \right), \end{aligned}$$

valid when Σ_0 and Σ_1 are symmetric positive definite so that all terms are finite.

For $d = 1$, set $\Sigma_0 = \sigma_0^2$ and $\Sigma_1 = \sigma_1^2$. Then $\text{tr}(\Sigma_1^{-1} \Sigma_0) = \sigma_0^2 / \sigma_1^2$, $\log \det$ becomes $\log(\sigma_1^2 / \sigma_0^2)$, and the quadratic form becomes $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) / \sigma_1^2$. Therefore

$$D_{KL}(\mathcal{N}(\boldsymbol{\mu}_0, \sigma_0^2) \| \mathcal{N}(\boldsymbol{\mu}_1, \sigma_1^2)) = \frac{1}{2} \left(\frac{\sigma_0^2}{\sigma_1^2} + \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)}{\sigma_1^2} - 1 + \log \frac{\sigma_1^2}{\sigma_0^2} \right).$$

Problem 2: Autoregressive (AR) Models of Order p

1,2. Let $\{y_t\}_{t=1}^n$ follow the AR(p) model

$$y_t = \sum_{j=1}^p \phi_j y_{t-j} + \varepsilon_t, \quad \varepsilon_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2).$$

Define one-step-ahead innovations

$$e_t(\phi) = y_t - \sum_{j=1}^p \phi_j y_{t-j}, \quad t = p+1, \dots, n,$$

where $\phi = (\phi_1, \dots, \phi_p)^\top$. Conditional on the initial lag vector (y_1, \dots, y_p) , the variables $e_t(\phi)$ are independent $\mathcal{N}(0, \sigma^2)$. Hence the conditional likelihood is

$$L_c(\phi, \sigma^2 | y_{1:n}) = \prod_{t=p+1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{e_t(\phi)^2}{2\sigma^2}\right) = (2\pi\sigma^2)^{-\frac{n-p}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{t=p+1}^n e_t(\phi)^2\right).$$

Therefore the conditional log-likelihood is

$$\begin{aligned} \ell_c(\phi, \sigma^2) &= -\frac{n-p}{2} \log(2\pi) - \frac{n-p}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} S(\phi), \\ S(\phi) &= \sum_{t=p+1}^n \left(y_t - \sum_{j=1}^p \phi_j y_{t-j} \right)^2. \end{aligned}$$

It is convenient to write this in matrix form. Let

$$\mathbf{y} = \begin{pmatrix} y_{p+1} \\ y_{p+2} \\ \vdots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} y_p & y_{p-1} & \cdots & y_1 \\ y_{p+1} & y_p & \cdots & y_2 \\ \vdots & \vdots & \ddots & \vdots \\ y_{n-1} & y_{n-2} & \cdots & y_{n-p} \end{pmatrix}, \quad \phi = \begin{pmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_p \end{pmatrix}.$$

Then $\mathbf{e}(\phi) = \mathbf{y} - X\phi$ and

$$S(\phi) = (\mathbf{y} - X\phi)^\top (\mathbf{y} - X\phi).$$

The log-likelihood becomes

$$\ell_c(\phi, \sigma^2) = -\frac{n-p}{2} \log(2\pi) - \frac{n-p}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - X\phi)^\top (\mathbf{y} - X\phi).$$

Maximization with respect to ϕ . For fixed σ^2 , maximizing ℓ_c in ϕ is equivalent to minimizing $S(\phi)$. Expand

$$S(\phi) = \mathbf{y}^\top \mathbf{y} - 2\phi^\top X^\top \mathbf{y} + \phi^\top X^\top X \phi.$$

Differentiate and set the gradient to zero:

$$\nabla_\phi S(\phi) = -2X^\top \mathbf{y} + 2X^\top X \phi = \mathbf{0} \iff X^\top X \phi = X^\top \mathbf{y}.$$

When $X^\top X$ is invertible, the unique solution of the normal equations is

$$\hat{\phi} = (X^\top X)^{-1} X^\top \mathbf{y}.$$

Equivalently, compute the log-likelihood gradient directly:

$$\nabla_{\phi} \ell_c(\phi, \sigma^2) = -\frac{1}{2\sigma^2} \nabla_{\phi} S(\phi) = -\frac{1}{\sigma^2} (X^\top X \phi - X^\top \mathbf{y}),$$

so the stationary condition is the same. The Hessian with respect to ϕ is

$$\nabla_{\phi}^2 \ell_c(\phi, \sigma^2) = -\frac{1}{\sigma^2} X^\top X,$$

which is negative definite when $X^\top X$ is positive definite. Hence $\hat{\phi}$ is the unique maximizer in ϕ for fixed σ^2 .

Maximization with respect to σ^2 . For fixed ϕ , differentiate ℓ_c with respect to σ^2 :

$$\frac{\partial}{\partial \sigma^2} \ell_c(\phi, \sigma^2) = -\frac{n-p}{2} \frac{1}{\sigma^2} + \frac{1}{2} \frac{S(\phi)}{(\sigma^2)^2}.$$

Set to zero and solve for σ^2 :

$$-(n-p)\sigma^2 + S(\phi) = 0 \implies \hat{\sigma}^2(\phi) = \frac{1}{n-p} S(\phi).$$

To verify a maximum, compute the second derivative

$$\frac{\partial^2}{\partial (\sigma^2)^2} \ell_c(\phi, \sigma^2) = \frac{n-p}{2} \frac{1}{(\sigma^2)^2} - \frac{S(\phi)}{(\sigma^2)^3},$$

which, evaluated at $\sigma^2 = \hat{\sigma}^2(\phi) = S(\phi)/(n-p)$, equals

$$-\frac{n-p}{2} \frac{1}{(\sigma^2)^2} < 0.$$

Thus $\hat{\sigma}^2(\phi)$ maximizes ℓ_c in σ^2 . Combining with $\hat{\phi}$ above and substituting $\phi = \hat{\phi}$ yields the joint maximizers

$$\hat{\phi} = (X^\top X)^{-1} X^\top \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n-p} (\mathbf{y} - X \hat{\phi})^\top (\mathbf{y} - X \hat{\phi}).$$

3. Let $\mathcal{F}_{t-1} = \sigma(y_{t-1}, y_{t-2}, \dots)$. The model is

$$y_t = \phi y_{t-1} + \varepsilon_t, \quad \varepsilon_t \text{ independent of } \mathcal{F}_{t-1}, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2).$$

For any Borel set $B \subset \mathbb{R}$,

$$\mathbb{P}(y_t \in B \mid \mathcal{F}_{t-1}) = \mathbb{P}(\phi y_{t-1} + \varepsilon_t \in B \mid \mathcal{F}_{t-1}) = \mathbb{P}(\varepsilon_t \in B - \phi y_{t-1}) = \mathbb{P}(y_t \in B \mid y_{t-1}),$$

since conditioning on \mathcal{F}_{t-1} fixes y_{t-1} and ε_t is independent of \mathcal{F}_{t-1} . Thus

$$\mathbb{P}(y_t \in B \mid y_{t-1}, y_{t-2}, \dots) = \mathbb{P}(y_t \in B \mid y_{t-1}),$$

so $\{y_t\}$ is a time-homogeneous Markov process of order 1.

(a) Conditional on $y_{t-1} = x$,

$$(y_t \mid y_{t-1} = x) = \phi x + \varepsilon_t \sim \mathcal{N}(\phi x, \sigma^2).$$

Hence, for $y \in \mathbb{R}$,

$$p(y_t = y \mid y_{t-1} = x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \phi x)^2}{2\sigma^2}\right),$$

which depends on t only through the constant parameters ϕ, σ^2 , so the chain is time-homogeneous.

(b) Iterate the recursion:

$$y_t = \phi^t y_0 + \sum_{k=0}^{t-1} \phi^k \varepsilon_{t-k}.$$

If $|\phi| < 1$, then $\phi^t y_0 \rightarrow 0$ almost surely. The sum is a linear combination of independent zero-mean Gaussians, hence Gaussian with mean 0 and variance

$$\text{Var}(y_t) = \sum_{k=0}^{t-1} \phi^{2k} \sigma^2 = \sigma^2 \sum_{k=0}^{t-1} \phi^{2k} = \sigma^2 \frac{1 - \phi^{2t}}{1 - \phi^2} \xrightarrow{t \rightarrow \infty} \frac{\sigma^2}{1 - \phi^2}.$$

Therefore the unique stationary distribution is

$$\pi(y) = \mathcal{N}\left(0, \frac{\sigma^2}{1 - \phi^2}\right),$$

with mean 0 and variance $\sigma^2 / (1 - \phi^2)$.

Problem 3: Nonlinear and Multi-Step Autoregressive Models

1. • Conditioning on the initial window $x_{1:p}$, the likelihood factors as

$$p_\theta(x_{p+1:T} \mid x_{1:p}) = \prod_{t=p+1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_t - f_\theta(x_{t-1:t-p}))^2}{2\sigma^2}\right).$$

Hence the negative log likelihood is

$$\mathcal{L}(\theta) = -\log p_\theta(x_{p+1:T} \mid x_{1:p}) = \frac{T-p}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{t=p+1}^T (x_t - f_\theta(x_{t-1:t-p}))^2.$$

Define the per step residual $r_t = f_\theta(x_{t-1:t-p}) - x_t$. Then

$$\mathcal{L}(\theta) = \frac{T-p}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{t=p+1}^T r_t^2.$$

- For a windowed feedforward f_θ that takes $(x_{t-1}, \dots, x_{t-p})$ as inputs, the gradient is a sum of per step contributions with shared weights

$$\nabla_\theta \mathcal{L}(\theta) = \frac{1}{2\sigma^2} \sum_{t=p+1}^T \frac{\partial r_t^2}{\partial r_t} \cdot \frac{\partial r_t}{\partial \theta} = \frac{1}{\sigma^2} \sum_{t=p+1}^T r_t \nabla_\theta f_\theta(x_{t-1:t-p}).$$

Computationally, one runs a forward pass to compute all $f_\theta(x_{t-1:t-p})$, accumulates the residuals r_t , and backpropagates through the network of f_θ at each time step, summing the parameter gradients because the same θ is reused.

- When training with teacher forcing, each $f_\theta(x_{t-1:t-p})$ consumes ground truth inputs, so there is no dependency of x_t on earlier model outputs within the loss. Gradients do not have to traverse across time through the data sequence; they only traverse within the network that defines f_θ . The only coupling across time is weight sharing, which results in summation of per step gradients.

When f_θ is recurrent or when training objectives unroll the model over its own predictions for multi step horizons, gradients propagate through time via products of Jacobians such as $J_{g,h}$. For long sequences or large effective horizons this can cause vanishing or exploding gradients depending on the spectral norms of these Jacobians. If the norms are mostly below one, products tend to zero and gradients vanish, hampering learning of long range dependencies. If they are above one, products can blow up and gradients explode. Common mitigations include residual or gating architectures, gradient clipping, careful initialization, normalization, and contracting regularization that keeps the state update near a contraction. For large p with a purely feedforward f_θ , the input dimension grows but the gradient path length in time does not, so the main concern shifts to optimization conditioning and input scaling rather than temporal vanishing or exploding.

2. • When the order p is very large, the nonlinear function f_θ receives a long input vector $\mathbf{x}_{t-1:t-p} = (x_{t-1}, \dots, x_{t-p})$. A fully connected multilayer perceptron (MLP) whose first layer directly connects to all p lags has $O(p)$ inputs and quickly becomes expensive and prone to overfitting as p grows. In this regime it is natural to exploit structure in time.

One important design is a one dimensional convolutional network over the lag axis. A 1D CNN with kernel size k and shared weights across time can model local temporal patterns with $O(k)$ parameters per filter instead of $O(p)$. Stacking layers increases the effective receptive field. With dilated convolutions the receptive field can grow exponentially in depth, so a deep dilated CNN (as in WaveNet type architectures) can cover a very large context with logarithmic depth. This yields

$$f_\theta(x_{t-1}, \dots, x_{t-p}) = \text{Conv1D}_\theta(x_{t-p:t-1}),$$

where the convolution is causal and may include residual connections and normalization to improve training stability.

Self attention based architectures can also be used when p is large. A transformer style block takes all lags x_{t-1}, \dots, x_{t-p} as a sequence and uses attention to adaptively focus on a subset of relevant time steps. This can capture long range dependencies and interactions that are hard for local convolutions. However, the naive attention

cost is $O(p^2)$ and memory usage also scales quadratically in p , which becomes prohibitive when p is very large. Practical use requires sparse or low rank attention, truncation of context, or streaming approximations.

Recurrent neural networks (RNNs) such as LSTMs or GRUs can be used to summarize the long history into a hidden state \mathbf{h}_{t-1} , and one can set

$$f_\theta(x_{t-1}, \dots, x_{t-p}) = g_\theta(\mathbf{h}_{t-1}), \quad \mathbf{h}_{t-1} = \text{RNN}_\theta(\mathbf{h}_{t-p-1}, x_{t-p:t-1}).$$

This avoids a p dimensional input layer, but gradients through long sequences may vanish or explode, and training cost scales with sequence length. Gated units and techniques like truncated backpropagation, gradient clipping, and normalization help but there remain practical limits on the effective horizon.

In summary, for $p \gg 1$ it is usually preferable to choose architectures with

- parameter sharing across time (convolutions, RNNs, attention),
- controlled growth of receptive field (dilated convolutions or attention),
- good inductive bias for temporal structure (causal convolutions or causal attention).

Fully connected MLPs with dense connections to all lags become inefficient and unstable as p grows.

- When p is small, for example $p = 1, 2, 3$, the input dimension to f_θ is low and the model does not need to handle very long temporal dependencies. In this regime simple architectures are often preferable.

A shallow MLP with a few hidden layers and nonlinearities is a natural choice. The mapping

$$f_\theta(x_{t-1}, \dots, x_{t-p}) = \text{MLP}_\theta(x_{t-1}, \dots, x_{t-p})$$

has

- good approximation power for smooth nonlinear relationships between a small number of lags,
- low computational cost per time step, since the first layer has only p inputs,
- simple and stable training behavior, because gradients do not need to propagate through time and the network is relatively shallow.

For $p = 1$, even a single hidden layer with a modest number of units is often sufficient, and adding recurrent or attention mechanisms usually brings little benefit while increasing complexity and risk of overfitting.

Simple recurrent models like GRUs or LSTMs can still be used for small p , but they are typically unnecessary compared to MLPs unless one wants to process variable length contexts or share a single recurrent state across many different predictive tasks. Transformers and attention mechanisms are generally excessive in this setting: the sequence length seen by f_θ is tiny, so the advantages of flexible long range attention do not outweigh their computational and implementation overhead.

3. • Define deterministic point predictions as conditional means. First step:

$$\hat{x}_{t+1} = f_\theta(x_t, x_{t-1}, \dots, x_{t-p+1}).$$

For later steps we feed back previous predictions. For $j \geq 2$,

$$\hat{x}_{t+j} = f_\theta(x_{t+j-1}^*, x_{t+j-2}^*, \dots, x_{t+p}^*),$$

where

$$x_s^* = \begin{cases} x_s & s \leq t, \\ \hat{x}_s & s > t. \end{cases}$$

Then $\hat{x}_{t+1}, \hat{x}_{t+2}, \dots, \hat{x}_{t+k}$ are defined recursively for any k .

- Assume the true process is

$$x_{t+1} = f_\star(x_t, \dots, x_{t-p+1}) + \varepsilon_{t+1}, \quad \mathbb{E}[\varepsilon_{t+1}] = 0.$$

Define the true state vector

$$s_t = \begin{pmatrix} x_t \\ x_{t-1} \\ \vdots \\ x_{t-p+1} \end{pmatrix} \in \mathbb{R}^p,$$

and the model state update map

$$F_\theta(s_t) = \begin{pmatrix} f_\theta(x_t, \dots, x_{t-p+1}) \\ x_t \\ \vdots \\ x_{t-p+2} \end{pmatrix}.$$

Similarly the true update is

$$F_\star(s_t) = \begin{pmatrix} f_\star(x_t, \dots, x_{t-p+1}) \\ x_t \\ \vdots \\ x_{t-p+2} \end{pmatrix}, \quad \eta_{t+1} = \begin{pmatrix} \varepsilon_{t+1} \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Then $s_{t+1} = F_\star(s_t) + \eta_{t+1}$. The free running model prediction obeys

$$\hat{s}_{t+1} = F_\theta(\hat{s}_t), \quad \hat{s}_t = \begin{pmatrix} x_t \\ x_{t-1} \\ \vdots \\ x_{t-p+1} \end{pmatrix}.$$

Define the state error $e_t = \hat{s}_t - s_t$. Then

$$\begin{aligned} e_{t+1} &= \hat{s}_{t+1} - s_{t+1} = F_\theta(\hat{s}_t) - F_\star(s_t) - \eta_{t+1} \\ &= F_\theta(\hat{s}_t) - F_\theta(s_t) + \underbrace{F_\theta(s_t) - F_\star(s_t)}_{\text{model bias}} - \eta_{t+1}. \end{aligned}$$

Assume F_θ is Lipschitz with constant L in some norm $\|\cdot\|$, and that the bias is bounded by $\delta \geq 0$,

$$\|F_\theta(u) - F_\theta(v)\| \leq L \|u - v\|, \quad \|F_\theta(s_t) - F_\star(s_t)\| \leq \delta.$$

Then

$$\|e_{t+1}\| \leq L \|e_t\| + \delta + \|\eta_{t+1}\|.$$

Iterating this recursion gives, for any $k \geq 1$,

$$\|e_{t+k}\| \leq L^k \|e_t\| + \sum_{j=0}^{k-1} L^j (\delta + \|\eta_{t+k-j}\|).$$

Even if $\|e_t\| = 0$ initially, the expectation of $\|e_{t+k}\|$ grows with k as a geometric sum of model bias and noise terms. Thus multi step prediction errors are generally much larger than one step errors.

One way to reduce error accumulation is to train for multi step prediction. Fix a horizon K . For each starting time t , define free running predictions $\hat{x}_{t+1}, \dots, \hat{x}_{t+K}$ by the recursion above and use the loss

$$\mathcal{L}_K(\theta) = \sum_t \sum_{j=1}^K (x_{t+j} - \hat{x}_{t+j})^2.$$

Backpropagation through time on $\mathcal{L}_K(\theta)$ forces the model to be robust to its own prediction errors and better aligned with the distribution of inputs it will see during multi step forecasting.

Another way is scheduled sampling. At training time, instead of always feeding the true past x_s into the model, randomly mix true and predicted values. For each step s define

$$z_s = m_s x_s + (1 - m_s) \hat{x}_s,$$

where $m_s \sim \text{Bernoulli}(q)$ is an independent mask with probability $q \in [0, 1]$ of using the true value. The model then uses windows built from z_s rather than x_s . Early in training one chooses $q \approx 1$ and gradually decreases q so that the model increasingly trains under its own predictions. This reduces the train test mismatch and mitigates error accumulation in long horizon forecasts.

4. • Conditioning on the first p values $x_{1:p}$, the conditional likelihood of $x_{p+1:T}$ given $c_{1:T}$ factorizes as

$$\begin{aligned} p_\theta(x_{p+1:T} | x_{1:p}, c_{1:T}) &= \prod_{t=p+1}^T p_\theta(x_t | x_{<t}, c_{1:t}) \\ &= \prod_{t=p+1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_t - f_\theta(x_{t-1:t-p}; c_{1:t}))^2}{2\sigma^2}\right). \end{aligned}$$

The conditional log likelihood is

$$\begin{aligned} \log p_\theta(x_{p+1:T} | x_{1:p}, c_{1:T}) &= \sum_{t=p+1}^T \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (x_t - f_\theta(x_{t-1:t-p}; c_{1:t}))^2 \right] \\ &= -\frac{T-p}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=p+1}^T (x_t - f_\theta(x_{t-1:t-p}; c_{1:t}))^2. \end{aligned}$$

For fixed σ^2 , maximizing this with respect to θ is equivalent to minimizing

$$\mathcal{L}_{\text{cond}}(\theta) = \frac{1}{2\sigma^2} \sum_{t=p+1}^T (x_t - f_\theta(x_{t-1:t-p}; c_{1:t}))^2,$$

or any positive constant multiple of this mean squared error. With a dataset of sequences indexed by n , the full training objective is the sum over sequences

$$\mathcal{L}_{\text{train}}(\theta) = \sum_n \sum_{t=p+1}^{T_n} \left(x_t^{(n)} - f_\theta(x_{t-1:t-p}^{(n)}; c_{1:t}^{(n)}) \right)^2.$$

- Given a trained model f_θ , a conditioning sequence $c_{1:T}$, and some initial values $x_{1:p}$ (either observed or sampled from an initial distribution), we can generate a conditional sequence by forward simulation.

For $t = p + 1, \dots, T$ do:

$$x_t \sim \mathcal{N}(f_\theta(x_{t-1}, \dots, x_{t-p}; c_{1:t}), \sigma^2),$$

feeding the newly generated x_t back into the inputs for future steps. Using the mean instead of sampling (that is, setting $x_t = f_\theta(\cdot)$) yields a deterministic conditional forecast.

For example, in conditional speech generation, $c_{1:T}$ might be a sequence of text or linguistic features and $x_{1:T}$ an audio waveform. The model learns how the waveform evolves over time given the text. At generation time, for a new text sequence $c_{1:T}$, we sample or predict $x_{1:T}$ from the conditional model. Similarly, for motion generation, $c_{1:T}$ can be control commands or trajectory waypoints, and $x_{1:T}$ can be joint angles or positions.

- An unconditional autoregressive model for $x_{1:T}$ factorizes

$$p_\theta(x_{1:T}) = \prod_{t=1}^T p_\theta(x_t | x_{<t}),$$

so all dependence is internal to the sequence x . In the conditional case, the conditional distribution factorizes as

$$p_\theta(x_{1:T} | c_{1:T}) = \prod_{t=1}^T p_\theta(x_t | x_{<t}, c_{1:t}).$$

Here x_t depends both on its own past $x_{<t}$ and on the conditioning signal $c_{1:t}$. If we regard c as exogenous, the joint density factorizes as

$$p(x_{1:T}, c_{1:T}) = p(c_{1:T}) p_\theta(x_{1:T} | c_{1:T}),$$

while in the unconditional case there is no explicit conditioning sequence.

In terms of dependence structure, the unconditional AR model encodes

$$x_t \perp\!\!\!\perp x_{>t} | x_{<t},$$

while the conditional AR additionally yields

$$x_t \perp\!\!\!\perp x_{>t} \mid x_{<t}, c_{1:T},$$

with direct functional dependence on $c_{1:t}$ at each step. This allows the model to explain variation in x_t caused by the external sequence c , rather than forcing all structure to be captured solely by past x .

Consequently, unconditional AR models are suitable for modeling a single sequence in isolation, such as raw audio or financial returns with no external covariates. Conditional AR models are designed for tasks where the sequence x is driven or modulated by another sequence c , such as speech given text, motion given control signals, or sensor readings given planned actions. The conditional formulation can capture richer input–output relationships and enables controllable generation by varying the conditioning sequence.

- Let $x_{1:T} = (x_1, \dots, x_T)$ and $c_{1:T} = (c_1, \dots, c_T)$. Consider two conditional autoregressive models

$$p_\theta(x_t \mid x_{<t}, c_{1:t}), \quad q_\phi(x_t \mid x_{<t}, c_{1:t}),$$

with corresponding sequence distributions

$$p_\theta(x_{1:T} \mid c_{1:T}) = \prod_{t=1}^T p_\theta(x_t \mid x_{<t}, c_{1:t}), \quad q_\phi(x_{1:T} \mid c_{1:T}) = \prod_{t=1}^T q_\phi(x_t \mid x_{<t}, c_{1:t}).$$

Fix a conditioning sequence $c_{1:T}$. The Kullback–Leibler divergence between the two conditional distributions is

$$\begin{aligned} D_{KL}(p_\theta(\cdot \mid c_{1:T}) \| q_\phi(\cdot \mid c_{1:T})) &= \int p_\theta(x_{1:T} \mid c_{1:T}) \log \frac{p_\theta(x_{1:T} \mid c_{1:T})}{q_\phi(x_{1:T} \mid c_{1:T})} dx_{1:T} \\ &= \mathbb{E}_{p_\theta(x_{1:T} \mid c_{1:T})} [\log p_\theta(x_{1:T} \mid c_{1:T}) - \log q_\phi(x_{1:T} \mid c_{1:T})]. \end{aligned}$$

Using the autoregressive factorization,

$$\begin{aligned} \log p_\theta(x_{1:T} \mid c_{1:T}) &= \sum_{t=1}^T \log p_\theta(x_t \mid x_{<t}, c_{1:t}), \\ \log q_\phi(x_{1:T} \mid c_{1:T}) &= \sum_{t=1}^T \log q_\phi(x_t \mid x_{<t}, c_{1:t}), \end{aligned}$$

so

$$D_{KL}(p_\theta \| q_\phi) = \sum_{t=1}^T \mathbb{E}_{p_\theta(x_{1:T} \mid c_{1:T})} [\log p_\theta(x_t \mid x_{<t}, c_{1:t}) - \log q_\phi(x_t \mid x_{<t}, c_{1:t})].$$

In the Gaussian NAR case,

$$p_\theta(x_t \mid x_{<t}, c_{1:t}) = \mathcal{N}(x_t \mid f_\theta(x_{t-1:t-p}; c_{1:t}), \sigma_p^2),$$

$$q_\phi(x_t \mid x_{<t}, c_{1:t}) = \mathcal{N}(x_t \mid g_\phi(x_{t-1:t-p}; c_{1:t}), \sigma_q^2),$$

the inner term is the KL between two univariate Gaussians for fixed history $(x_{<t}, c_{1:t})$,

$$D_{KL}(p_\theta(\cdot | x_{<t}, c_{1:t}) \| q_\phi(\cdot | x_{<t}, c_{1:t})) = \frac{1}{2} \left(\log \frac{\sigma_q^2}{\sigma_p^2} + \frac{\sigma_p^2 + (f_\theta - g_\phi)^2}{\sigma_q^2} - 1 \right),$$

where $f_\theta = f_\theta(x_{t-1:t-p}; c_{1:t})$ and $g_\phi = g_\phi(x_{t-1:t-p}; c_{1:t})$. The full sequence level KL is then

$$D_{KL}(p_\theta \| q_\phi) = \sum_{t=1}^T \mathbb{E}_{p_\theta(x_{1:T} | c_{1:T})} [D_{KL}(p_\theta(\cdot | x_{<t}, c_{1:t}) \| q_\phi(\cdot | x_{<t}, c_{1:t}))].$$

Computational challenges. Formally, the expectation is over the joint $p_\theta(x_{1:T} | c_{1:T})$ on \mathbb{R}^T . Even though this distribution factorizes autoregressively, the histories $(x_{<t}, c_{1:t})$ enter nonlinearly through f_θ and g_ϕ , so there is no closed form for the expectation in general. Direct evaluation would require integrating over all $x_{1:T}$, which is intractable for high dimensional T . Dynamic programming does not help because the log ratio involves both models, and only one of them defines the sampling distribution p_θ . For neural network f_θ, g_ϕ , the integrand does not admit analytic marginalization.

Monte Carlo approximation. A practical approach is Monte Carlo estimation using ancestral sampling from p_θ . For a fixed conditioning sequence $c_{1:T}$:

1. Draw N independent trajectories $\mathbf{x}_{1:T}^{(n)} \sim p_\theta(x_{1:T} | c_{1:T})$ by forward simulation

$$x_t^{(n)} \sim p_\theta(\cdot | x_{<t}^{(n)}, c_{1:t}), \quad t = 1, \dots, T.$$

2. For each sample compute

$$\log p_\theta(\mathbf{x}_{1:T}^{(n)} | c_{1:T}) = \sum_{t=1}^T \log p_\theta(x_t^{(n)} | x_{<t}^{(n)}, c_{1:t}),$$

$$\log q_\phi(\mathbf{x}_{1:T}^{(n)} | c_{1:T}) = \sum_{t=1}^T \log q_\phi(x_t^{(n)} | x_{<t}^{(n)}, c_{1:t}).$$

3. Approximate the KL divergence by the sample mean

$$D_{KL}(p_\theta \| q_\phi) \approx \frac{1}{N} \sum_{n=1}^N \left[\log p_\theta(\mathbf{x}_{1:T}^{(n)} | c_{1:T}) - \log q_\phi(\mathbf{x}_{1:T}^{(n)} | c_{1:T}) \right].$$

This estimator converges to the true KL as $N \rightarrow \infty$ by the law of large numbers. For random conditioning sequences, one can also average this estimate over a dataset of $c_{1:T}$ drawn from the empirical distribution, yielding an approximation of $\mathbb{E}_c [D_{KL}(p_\theta(\cdot | c) \| q_\phi(\cdot | c))]$ in high dimensional autoregressive settings.

Problem 4: Real NADE Parameters

Let $\mathbf{h}_i \in \mathbb{R}^d$ be the hidden representation computed from $\mathbf{x}_{<i}$. We model the scalar conditional as a C component Gaussian mixture

$$p(x_i | \mathbf{x}_{<i}) = \sum_{c=1}^C \pi_i^c \mathcal{N}(x_i | \mu_i^c, \exp(s_i^c)),$$

where π_i^c are mixture weights, μ_i^c are means, and s_i^c are log variances so that the variance of component c is $\exp(s_i^c)$.

Matrix parameterization from \mathbf{h}_i . Introduce parameter matrices and biases

$$W_i^{(\pi)} \in \mathbb{R}^{C \times d}, \quad \mathbf{a}_i^{(\pi)} \in \mathbb{R}^C, \quad V_i \in \mathbb{R}^{C \times d}, \quad \mathbf{b}_i \in \mathbb{R}^C, \quad U_i \in \mathbb{R}^{C \times d}, \quad \mathbf{d}_i \in \mathbb{R}^C.$$

Compute componentwise logits, means, and log variances as

$$\mathbf{z}_i = W_i^{(\pi)} \mathbf{h}_i + \mathbf{a}_i^{(\pi)}, \quad \boldsymbol{\mu}_i = V_i \mathbf{h}_i + \mathbf{b}_i, \quad \mathbf{s}_i = U_i \mathbf{h}_i + \mathbf{d}_i.$$

Let z_i^c, μ_i^c, s_i^c denote the c th entries of $\mathbf{z}_i, \boldsymbol{\mu}_i, \mathbf{s}_i$. Map logits to mixture weights with a softmax so that weights are nonnegative and sum to one

$$\pi_i^c = \frac{\exp(z_i^c)}{\sum_{k=1}^C \exp(z_i^k)}, \quad c = 1, \dots, C.$$

This choice ensures $\sum_{c=1}^C \pi_i^c = 1$ automatically and permits unconstrained optimization of z_i^c .

Equivalent per component vector form. For each c introduce $\mathbf{w}_i^c, \mathbf{v}_i^c, \mathbf{u}_i^c \in \mathbb{R}^d$ and scalars a_i^c, b_i^c, d_i^c and set

$$z_i^c = (\mathbf{w}_i^c)^\top \mathbf{h}_i + a_i^c, \quad \mu_i^c = (\mathbf{v}_i^c)^\top \mathbf{h}_i + b_i^c, \quad s_i^c = (\mathbf{u}_i^c)^\top \mathbf{h}_i + d_i^c, \quad \pi_i^c = \frac{\exp(z_i^c)}{\sum_{k=1}^C \exp(z_i^k)}.$$

This is exactly the same as the matrix form with rows of the matrices collected into the per component vectors.

Parameter count for a single conditional $p(x_i | \mathbf{x}_{<i})$. The mapping $\mathbf{h}_i \mapsto \{\pi_i^c, \mu_i^c, s_i^c\}_{c=1}^C$ uses

$$\underbrace{C(d+1)}_{\text{logits } \mathbf{z}_i} + \underbrace{C(d+1)}_{\text{means } \boldsymbol{\mu}_i} + \underbrace{C(d+1)}_{\text{log variances } \mathbf{s}_i} = 3C(d+1)$$

learned parameters for that single conditional. Although the simplex of $\boldsymbol{\pi}_i$ has $C - 1$ degrees of freedom, the softmax parameterization with C logits is standard, smooth, and avoids constraints. The upstream parameters that produce \mathbf{h}_i in Real NADE are shared across i and are not included in this per conditional count.

Notes. If desired, one can tie parameters across c or across i to reduce the count, or add temperature scaling to \mathbf{z}_i for sharper or flatter mixtures. For multivariate x_i , the same scheme applies but each component would output vector means and either diagonal or full covariance parameters.

Problem 5: Autoregressive Graph Neural Networks for Sequential Data

1. Let \mathcal{V} be the set of nodes and for each $v \in \mathcal{V}$ let x_v^t be the observation at time t . Define hidden states h_v^t updated by a GNN as

$$h_v^t = \text{GNN}(h_v^{t-1}, \{h_u^{t-1} : u \in N(v)\}), \quad t = 2, \dots, T,$$

with some initialization for h_v^1 . Assume an emission model $p(x_v^t | h_v^t)$, and that, conditional on $\{h_v^t\}_{v \in \mathcal{V}}$, the observations at time t are independent across nodes:

$$p(x_t | h_t) = \prod_{v \in \mathcal{V}} p(x_v^t | h_v^t), \quad x_t = (x_v^t)_{v \in \mathcal{V}}, \quad h_t = (h_v^t)_{v \in \mathcal{V}}.$$

First factorize over time using the chain rule:

$$p(x_{1:T}) = p(x_1) \prod_{t=2}^T p(x_t | x_{<t}), \quad x_{1:T} = (x_1, \dots, x_T).$$

For a fixed $t \geq 2$, by conditional independence of nodes given h_t ,

$$p(x_t | h_t) = \prod_{v \in \mathcal{V}} p(x_v^t | h_v^t).$$

The hidden state h_v^t is a deterministic function of the previous hidden states of v and its neighbors:

$$h_v^t = F_v(h_v^{t-1}, \{h_u^{t-1} : u \in N(v)\}).$$

Unrolling this recursion shows that h_v^t is in fact a deterministic function of the past observations of v and its neighbors. Thus there exists a function H_v such that

$$h_v^t = H_v(x_v^{<t}, x_{N(v)}^{<t}),$$

where $x_v^{<t} = (x_v^1, \dots, x_v^{t-1})$ and $x_{N(v)}^{<t} = \{x_u^{1:t-1} : u \in N(v)\}$. Therefore

$$p(x_v^t | x_{<t}) = p(x_v^t | h_v^t) = p(x_v^t | x_v^{<t}, x_{N(v)}^{<t}),$$

and, given the past, the nodewise observations at time t are conditionally independent:

$$p(x_t | x_{<t}) = \prod_{v \in \mathcal{V}} p(x_v^t | x_{<t}) = \prod_{v \in \mathcal{V}} p(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

Substitute this into the temporal factorization to obtain

$$p(x_{1:T}) = p(x_1) \prod_{t=2}^T \prod_{v \in \mathcal{V}} p(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

If we absorb the initial time $t = 1$ into the same notation by defining suitable initial factors, we can write the joint likelihood factorization as

$$p(x_{1:T}) = \prod_{t=1}^T \prod_{v \in \mathcal{V}} p(x_v^t | x_v^{<t}, x_{N(v)}^{<t}),$$

which is the desired result.

2. Let \mathcal{V} be the node set. At time t , node v has scalar value x_v^t and binary mask $m_v^t \in \{0, 1\}$. The observed and missing sets are

$$O_t = \{v : m_v^t = 1\}, \quad M_t = \{v : m_v^t = 0\}.$$

Assume an autoregressive GNN model with factorization

$$p_\theta(\mathbf{x}_{1:T}) = \prod_{t=1}^T \prod_{v \in \mathcal{V}} p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}),$$

where $\mathbf{x}_t = (x_v^t)_{v \in \mathcal{V}}$.

- When only $\mathbf{x}_{O_t}^t$ is observed, the conditional log likelihood contribution at time t is

$$\log p_\theta(\mathbf{x}_{O_t}^t | \mathbf{x}_{<t}) = \sum_{v \in O_t} \log p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

Equivalently, using masks,

$$\log p_\theta(\mathbf{x}_{\text{obs}}) = \sum_{t=1}^T \sum_{v \in \mathcal{V}} m_v^t \log p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

The training objective is to maximize this log likelihood over θ , or minimize the masked negative log likelihood

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \sum_{v \in \mathcal{V}} m_v^t \log p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

To encourage robustness to missing inputs, we can:

- Inject artificial masks during training: randomly set some $m_v^t = 0$ even when the value is available, and exclude those terms from the loss via the mask.
- Feed the mask into the GNN as an additional feature, for example updating node states using inputs of the form $(m_v^{t-1} x_v^{t-1}, m_v^{t-1})$ and $(m_u^{t-1} x_u^{t-1}, m_u^{t-1})$ for neighbors $u \in N(v)$.

This forces the model to learn to make good predictions even when some neighbors or self values are missing or unreliable.

- At a given time t , suppose only $\mathbf{x}_{O_t}^t$ are observed, but the past $\mathbf{x}_{<t}$ is (possibly partially) known or has already been imputed. With the autoregressive factorization

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{<t}) = \prod_{v \in \mathcal{V}} p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}),$$

the missing values at time t are conditionally independent of each other and of the observed values at time t , given the past, so

$$p_\theta(\mathbf{x}_{M_t}^t | \mathbf{x}_{O_t}^t, \mathbf{x}_{<t}) = \prod_{v \in M_t} p_\theta(x_v^t | x_v^{<t}, x_{N(v)}^{<t}).$$

Thus each missing node $v \in M_t$ can be imputed independently from its predictive distribution. For example, in a Gaussian output model

$$p_{\theta} \left(x_v^t \mid x_v^{<t}, x_{N(v)}^{<t} \right) = \mathcal{N} \left(x_v^t \mid \mu_v^t, (\sigma_v^t)^2 \right),$$

we can either:

$$\hat{x}_v^t = \mu_v^t \quad (\text{conditional mean imputation}), \quad \text{or} \quad x_v^t \sim \mathcal{N} \left(\mu_v^t, (\sigma_v^t)^2 \right) \quad (\text{sampling}).$$

For sequences with missing values over many times, we proceed sequentially in t :

- At each time t , use the current history $\mathbf{x}_{<t}$ (containing true or previously imputed values) to compute the predictive distributions $p_{\theta} \left(x_v^t \mid x_v^{<t}, x_{N(v)}^{<t} \right)$.
- For $v \in O_t$ keep the observed x_v^t ; for $v \in M_t$ replace x_v^t by its mean or a sample from its predictive distribution.

In this way, the learned autoregressive GNN provides a principled mechanism to infer or sample missing node values using the conditional distributions it has learned from data.