

Lecture 2: Mathematical Preliminaries for ML

Shukai Gong

Not all required mathematical preliminaries are included in the sections above. It's just a reminder of some of my rusty math knowledge.

1 Norm

1. Vector Norms: A norm is a function $\|\cdot\| : \mathbb{C}^N \mapsto \mathbb{R}$ satisfying

1. $\|\mathbf{x}\| \geq 0$ and $\|\mathbf{x}\| = 0$ iff $\mathbf{x} = \mathbf{0}$
2. $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$ for any $\alpha \in \mathbb{C}$
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality)

Commonly-used Vector Norms:

- ℓ_0 -norm: $\|\mathbf{x}\|_0 = \sum_{i=1}^N \mathbb{I}(x_i \neq 0)$ (number of non-zero elements)
- ℓ_1 -norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^N |x_i|$, ℓ_2 -norm: $\|\mathbf{x}\|_2 = (\sum_{i=1}^N |x_i|^2)^{\frac{1}{2}}$, ℓ_p -norm: $\|\mathbf{x}\|_p = (\sum_{i=1}^N |x_i|^p)^{\frac{1}{p}}$
- ℓ_∞ -norm: $\|\mathbf{x}\|_\infty = \max_{i=1, \dots, N} |x_i|$
- Weighted norm: $\|\mathbf{x}\|_{\mathbf{W}} = \|\mathbf{W}\mathbf{x}\|$
- Frobenius norm: $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^M \sum_{j=1}^N |x_{ij}|^2}$

2. Induced Matrix Norms: $\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$

- $\|\mathbf{A}\|_1 = \max_i \|\mathbf{a}_i\|_1$, where \mathbf{a}_i is the i -th column of \mathbf{A}
- $\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^H \mathbf{A})}$, where $\rho(\cdot)$ is the spectral radius.
- $\|\mathbf{A}\|_\infty = \max_j \|\mathbf{A}_j\|_1$, where \mathbf{A}_j is the j -th row of \mathbf{A}

Some properties of induced matrix norms:

- Consistency: $\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\|\|\mathbf{x}\|$, $\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|$
- If $\mathbf{\Lambda}$ is a diagonal matrix, then $\|\mathbf{\Lambda}\|_p = \max_i |d_{ii}|$
- When $\mathbf{A} = \mathbf{a}$ is a vector, $\|\mathbf{A}\|_2 = \|\mathbf{a}\|_2$
- If $\mathbf{A} = \mathbf{u}\mathbf{v}^H$ is a rank-1 matrix, then $\|\mathbf{A}\|_2 = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$

2 Algebra

1. **Vector products:** given $\mathbf{a}, \mathbf{b} \in \mathbb{R}^N$

1. **Inner product:** $\mathbf{a}^\top \mathbf{b} = \mathbf{a} \cdot \mathbf{b} \langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^N a_i b_i$
2. **Outer product:** $\mathbf{a} \mathbf{b}^\top = \mathbf{a} \otimes \mathbf{b} = [a_i b_j] \in \mathbb{R}^{N \times N}$

2. **Matrix products:** given $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{M \times N}$

1. **Inner product:** $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B}) = \sum_{i=1}^M \sum_{j=1}^N a_{ij} b_{ij}$
2. **Outer product:** $\mathbf{A} \otimes \mathbf{B} = [a_{ij} b_{kl}] \in \mathbb{R}^{MN \times MN}$
3. **Kronecker product:** $\mathbf{A} \otimes \mathbf{B} = [a_{ij} b_{kl}] \in \mathbb{R}^{M^2 \times N^2}$

3. **Orthogonal vector and matrix:**

- **Orthogonal vector:**

- $\mathbf{a}, \mathbf{b} \in \mathbb{R}^N$ are orthogonal if $\mathbf{a}^\top \mathbf{b} = 0$ and $\|\mathbf{a}\|_2 \neq 0, \|\mathbf{b}\|_2 \neq 0$
- $\mathbf{a}, \mathbf{b} \in \mathbb{C}^N$ are orthogonal if $\mathbf{a}^H \mathbf{b} = \sum_{i=1}^N \bar{a}_i b_i = 0$ and $\|\mathbf{a}\|_2 \neq 0, \|\mathbf{b}\|_2 \neq 0$
- H represents Hermitian transpose for matrix (or vector) $\mathbf{A} \in \mathbb{C}^{M \times N}$

- **Orthonormal vector:** \mathbf{a}, \mathbf{b} are orthonormal if $\mathbf{a}^\top \mathbf{b} = 0$ and $\|\mathbf{a}\|_2 = \|\mathbf{b}\|_2 = 1$

- **Orthogonal matrix:** \mathbf{A} is orthogonal if $\mathbf{A}^\top \mathbf{A} = \mathbf{I}$, i.e. $\mathbf{A}^\top = \mathbf{A}^{-1}$

- Norm preserving: $\|\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ if \mathbf{A} is orthogonal

4. **Inverse Matrix:** A matrix which has an inverse is called **non-singular**, otherwise **singular**. \mathbf{A}^{-1} exists $\iff \det(\mathbf{A}) \neq 0$

- **Ill-conditioned matrix:** \mathbf{A} is close to being singular, that $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ is large.

- **Pseudo-inverse:** $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$, where \mathbf{A} could be non-square.

- It can be shown that $\mathbf{A}^+ \mathbf{A} = \mathbf{I}$ provided that $\mathbf{A}^\top \mathbf{A}$ is non-singular.

5. **Linear Systems:** $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{b} \in \mathbb{R}^M$

- \mathbf{A} is a **linear mapping:** $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$

- \mathbf{b} is a **linear combination** of columns of \mathbf{A} : $\mathbf{b} = \sum_{i=1}^N x_i \mathbf{a}_i$

There are 3 key tasks of linear systems

$$\underbrace{\mathbf{A}}_{\text{System}} \underbrace{\mathbf{x}}_{\text{Input}} = \underbrace{\mathbf{b}}_{\text{Output}}$$

- **Inverse Problem:** Given \mathbf{A} and \mathbf{b} , solve \mathbf{x} or $\min_{\mathbf{x}} d(\mathbf{A}\mathbf{x}, \mathbf{b})$

- **Modelling:** Given sets of \mathbf{b} 's (Denote \mathbf{B}) and \mathbf{x} 's (Denote \mathbf{X}), solve \mathbf{A} or $\min_{\mathbf{A}} d(\mathbf{A}\mathbf{X}, \mathbf{B})$

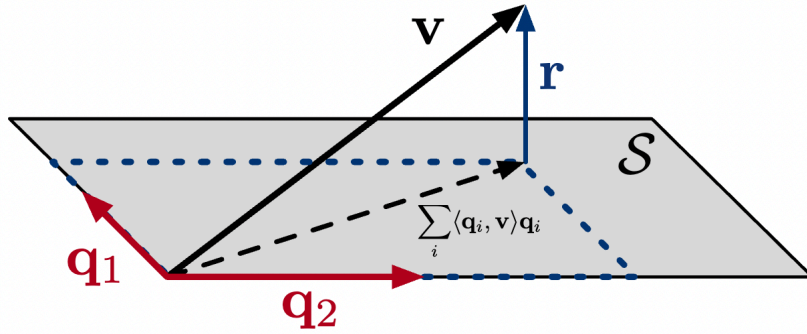
- **Factorization:** Given \mathbf{B} , solve the decomposition/factorization $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}, \mathbf{X}} d(\mathbf{A}\mathbf{X}, \mathbf{B})$

6. Range and Null Space:

- $\text{Range}(\mathbf{A}) = \text{Im}(\mathbf{A}) = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^N\}$, i.e. the column space of \mathbf{A}
- $\text{Null}(\mathbf{A}) = \{\mathbf{x} : \mathbf{Ax} = \mathbf{0}\}$

7. Components of a vector: Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ where $\mathbf{q}_1, \dots, \mathbf{q}_N$ is an orthonormal set in \mathbb{R}^m . Then for any $\mathbf{v} \in \mathbb{R}^m$, we have the **decomposition of \mathbf{v}** as shown below

$$\mathbf{v} = \sum_{i=1}^N \underbrace{\langle \mathbf{v}, \mathbf{q}_i \rangle}_{\in \mathcal{S}} \mathbf{q}_i + \underbrace{\mathbf{r}}_{\in \mathcal{S}^\perp}$$



- The residual of $\mathbf{v} \in \mathbb{R}^m$ w.r.t. the set $\mathbf{q}_1, \dots, \mathbf{q}_N$: $\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{v}, \mathbf{q}_i \rangle \mathbf{q}_i$

2.1 Eigenvalue Decomposition(EVD)

Eigenvalue Decomposition/Spectrum Decomposition

If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is **symmetric**, denote the orthonormal eigenvectors of \mathbf{A} as $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_n]$, and the corresponding eigenvalues as $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$, we can write

$$\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$$

The steps of EVD are as follows

1. **Find the eigenvalues:** By solving the characteristic equation $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$, we get $\lambda_1, \dots, \lambda_s$ where the sum of all algebraic multiplicity $m(\lambda_i)$ is n .
2. **Find the eigenvectors:** For each λ_i , solve the equation $(\mathbf{A} - \lambda_i \mathbf{I}) \mathbf{q}_i = \mathbf{0}$, we get $\mathbf{v}_{11}, \dots, \mathbf{v}_{1n_1}, \dots, \mathbf{v}_{s1}, \dots, \mathbf{v}_{sn_s}$ where the sum of all algebraic multiplicity $m(\lambda_i)$ is n .
3. **Orthogonalize the eigenvectors:** Use Gram-Schmidt process to orthogonalize $\mathbf{v}_{11}, \dots, \mathbf{v}_{1n_1}, \dots, \mathbf{v}_{s1}, \dots, \mathbf{v}_{sn_s}$ into $\mathbf{q}_{11}, \dots, \mathbf{q}_{1n_1}, \dots, \mathbf{q}_{s1}, \dots, \mathbf{q}_{sn_s}$
4. **Form the matrix \mathbf{Q} :** $\mathbf{Q} = [\mathbf{q}_{11}, \dots, \mathbf{q}_{1n_1}, \dots, \mathbf{q}_{s1}, \dots, \mathbf{q}_{sn_s}]$, we have $\mathbf{Q}^\top \mathbf{A} \mathbf{Q} = \text{diag}(\lambda_1, \dots, \lambda_n) = \mathbf{\Lambda} \Rightarrow \mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$

2.2 Singular Value Decomposition(SVD)

Singular Value Decomposition(SVD)

For any $\mathbf{A} \in \mathbb{R}^{m \times n}$, we can write

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$$

where $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$ are orthogonal matrices, and $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is a diagonal matrix with the **singular values of \mathbf{A}** σ_i on its diagonal, $\mathbf{u}_i, \mathbf{v}_i$ are the i -th columns of \mathbf{U} and \mathbf{V} . Only the first $r = \text{rank}(\mathbf{A})$ singular values are non-zero and by convention, they are ordered in non-increasing order: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_{\min\{m,n\}} = 0$.

Observe that the SVD factors provide eigendecomposition for $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{A}\mathbf{A}^\top$:

$$\begin{aligned} \mathbf{A}^\top \mathbf{A} &= (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top) = \mathbf{V}\mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V}(\mathbf{\Sigma}^\top \mathbf{\Sigma})\mathbf{V}^\top = \mathbf{V}\mathbf{\Lambda}_1\mathbf{V}^\top \\ \mathbf{A}\mathbf{A}^\top &= (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \mathbf{V}\mathbf{\Sigma}^\top \mathbf{U}^\top = \mathbf{U}(\mathbf{\Sigma}\mathbf{\Sigma}^\top)\mathbf{U}^\top \equiv \mathbf{U}\mathbf{\Lambda}_2\mathbf{U}^\top \end{aligned}$$

It follows immediately that the columns of \mathbf{V} are eigenvectors of $\mathbf{A}^\top \mathbf{A}$ and the columns of \mathbf{U} are eigenvectors of $\mathbf{A}\mathbf{A}^\top$.

The non-zero singular values of \mathbf{A} are the square roots of the non-zero eigenvalues of $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{A}\mathbf{A}^\top$.

The steps of SVD are as follows

1. **Find the orthogonal eigenvectors of $(\mathbf{A}^\top \mathbf{A})_{n \times n}$:** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$, where \mathbf{v}_i is the i -th orthogonal eigenvector of $\mathbf{A}^\top \mathbf{A}$.
2. **Find the orthogonal eigenvectors of $(\mathbf{A}\mathbf{A}^\top)_{m \times m}$:** $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$, where \mathbf{u}_i is the i -th orthogonal eigenvector of $\mathbf{A}\mathbf{A}^\top$.
3. **Form the singular value matrix $\mathbf{\Sigma}$:** $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \Rightarrow \mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma} \Rightarrow \mathbf{A}\mathbf{v}_i = \sigma_i \mathbf{u}_i \Rightarrow \sigma_i = \sqrt{\frac{\|\mathbf{A}\mathbf{v}_i\|_2}{\|\mathbf{v}_i\|_2}}$, and then we can form $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$

3 Matrix Calculus

1. Gradient:

- **Matrix Gradient:** Suppose that $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$, then the gradient of f w.r.t. $\mathbf{A} \in \mathbb{R}^{m \times n}$ is defined as

$$\nabla_{\mathbf{A}} f(\mathbf{A}) = \begin{bmatrix} \frac{\partial f(\mathbf{A})}{\partial a_{11}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{A})}{\partial a_{m1}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{mn}} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

- **Vector Gradient:** Suppose that $f : \mathbb{R}^n \mapsto \mathbb{R}$, then the gradient of f w.r.t. $\mathbf{x} \in \mathbb{R}^n$ is defined as

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

- Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}^m$, then the gradient of f w.r.t. $\mathbf{x} \in \mathbb{R}^n$ is defined as

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x})^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_m(\mathbf{x})^\top \end{bmatrix} = \mathbf{J}_f^\top(\mathbf{x}) \in \mathbb{R}^{m \times n}$$

2. Basic Facts about Matrix Derivatives:

1. $\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$
2. $\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$
3. $\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top$
4. $\frac{\partial \mathbf{a}^\top \mathbf{X}^\top \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^\top$
5. $\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^\top \mathbf{X}^\top \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^\top$

3. **Hessian Matrix:** Suppose that $f : \mathbb{R}^n \mapsto \mathbb{R}$, then the Hessian matrix of f w.r.t. $\mathbf{x} \in \mathbb{R}^n$ is defined as

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

4 Probability and Statistics

A typical ML scenario $(\mathcal{X}, d_{\mathcal{X}}, \mu_{\mathcal{X}})$ requires us to estimate $\mu_{\mathcal{X}}$ via a model \hat{p}_{θ} based on data $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{X}$.

4.1 Law of Large Numbers and Central Limit Theorem

1. **LLN explains why ML requires lots of data:** For $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{X}$

- WLLN: $\bar{X}_n \xrightarrow{P} \mu$
- SLLN: $\bar{X}_n \xrightarrow{a.s.} \mu$
- **Variance reduction:** $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \rightarrow 0$ as $n \rightarrow \infty$

2. **CLT provides ML with Gaussian Distribution:**

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1)$$

4.2 Method of Moments(MoM)

Suppose that we have a set of i.i.d. data $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{X}$ drawn from a distribution $P(\mathbf{X}|\boldsymbol{\theta})$ with l parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_l)$.

1. Compute the first l moments as functions of $\boldsymbol{\theta}$

$$\begin{cases} \mu_1 = E[X] = \int_{\mathcal{X}} \mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} = g_1(\theta_1, \dots, \theta_n) \\ \mu_2 = E[X^2] = \int_{\mathcal{X}} \mathbf{x}^2 P(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} = g_2(\theta_1, \dots, \theta_n) \\ \dots \\ \mu_l = E[X^l] = \int_{\mathcal{X}} \mathbf{x}^l P(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} = g_l(\theta_1, \dots, \theta_n) \end{cases}$$

2. Algebraically invert the linear system of l equations to solve for $\theta_1, \dots, \theta_l$ as functions of μ_1, \dots, μ_l

$$\begin{cases} \theta_1 = h_1(\mu_1, \dots, \mu_l) \\ \theta_2 = h_2(\mu_1, \dots, \mu_l) \\ \dots \\ \theta_l = h_l(\mu_1, \dots, \mu_l) \end{cases} \quad (*)$$

3. Insert **sample moments** $\hat{\mu}_k = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^k$ into (*) to obtain **MoM estimators** $\hat{\theta}_1, \dots, \hat{\theta}_l$

$$\begin{cases} \hat{\theta}_1 = h_1(\hat{\mu}_1, \dots, \hat{\mu}_l) \\ \hat{\theta}_2 = h_2(\hat{\mu}_1, \dots, \hat{\mu}_l) \\ \dots \\ \hat{\theta}_l = h_l(\hat{\mu}_1, \dots, \hat{\mu}_l) \end{cases}$$

Drawbacks of MoM

1. High computation load
2. Lack of extensibility
3. MoM estimators may not exist as the linear system of $\boldsymbol{\mu} = \mathbf{G}\boldsymbol{\theta}$ may have no solution \Rightarrow Can only be used to estimate simple distributions with few parameters

4.3 Maximum Likelihood Estimation

Suppose that we have a set of i.i.d. data $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{X}$, we assume that the samples are sampled from a distribution $P(\mathbf{x}|\theta)$ (a model with parameter θ).

Principle: Assume a **deterministic model** and learn the model via maximum likelihood estimation (MLE)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \prod_{i=1}^N \log P(\mathbf{x}_i|\theta) = \arg \max_{\theta} \sum_{i=1}^N \log P(\mathbf{x}_i|\theta)$$

- **Pros:** more efficient in general, avoid the design of prior.
- **Cons:** non-robust to sparse data, not easy to quantify the uncertainty of the estimation (doable, but not efficient).

4.4 Bayesian Estimation

Principle: Assume a **probabilistic model** and the model θ yields a prior distribution.

According to the Bayes' theorem, we have

$$\underbrace{P(\theta|\mathbf{X})}_{\text{Posterior}} = \frac{P(\mathbf{X}|\theta)P(\theta)}{P(\mathbf{X})} \propto \underbrace{P(\mathbf{X}|\theta)}_{\text{Likelihood } \theta} \underbrace{P(\theta)}_{\text{prior } \theta}$$

- The influence of prior decays with the increase of the number of samples.

MAP estimation: $\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|\mathbf{X}) = \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta)$

- **Pros:** prior makes it (relatively) robust to sparse data, quantify the uncertainty of the estimation (obtain the distribution of θ)
- **Cons:** require sophisticated design of prior, time-consuming in general.