

Maximum Likelihood Estimation (MLE)

Setup: Given some data $\mathcal{D} := \{x_1, x_2, \dots, x_n\}$, $x_i \in \mathbb{R}^d$

Assume a family of distributions/models $p_\theta(x)$, $\theta \in \Theta$

i.e., $x_i \sim p_\theta(x)$ for some θ .

Goal: Estimate the true value of θ that best explains the observed data.

Definition: $\hat{\theta}_{MLE}$ is a maximum likelihood estimate if:

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \underbrace{p(\mathcal{D}|\theta)}_{\text{likelihood}}, \quad p(\mathcal{D}|\hat{\theta}_{MLE}) = \max_{\theta} p(\mathcal{D}|\theta)$$

$$\textcircled{*} \quad p(\mathcal{D}|\theta) = p(x_1, x_2, \dots, x_n | \theta) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n p(x_i | \theta) = \prod_{i=1}^n p(X=x_i | \theta)$$

Remarks: i.) The MLE might not be unique.

ii.) The MLE may fail to exist (the maximum likelihood may not be achieved for $\theta \in \Theta$)

Pros: i.) Usually easy to compute and often

is interpretable (e.g. the mean or r.v. is the sample mean)

ii.) Nice asymptotic properties:

- Consistent: as $n \rightarrow +\infty$ the MLE converge to the true θ in probability.

... .. distribution

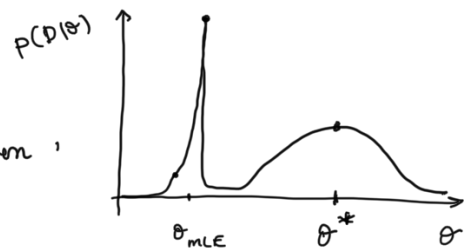
- asymptotically normal, as $n \rightarrow \infty$ it converges to a normal.
- efficient, i.e. they have the lowest asymptotic variance.
- invariant to re-parametrization:
 $\theta_{MLE}, \forall g: g(\theta_{MLE})$ is an MLE for $g(\theta)$.

Cons:

- i) MLE provides a point estimate (no representation of uncertainty)
- Ideally, we'd like to compute the posterior distribution over θ : $p(\theta|D) \stackrel{\text{Bayes}}{\propto} \underbrace{p(D|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}} \quad (\text{Bayesian approach})$
- Frequentist approach: "Train" an ensemble of models over many initialization $\theta_0 \sim p(\theta_0)$.

ii.) Lack of robustness:

e.g. thing like this may happen,

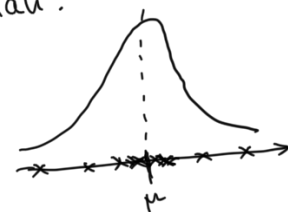


iii) MLE are prone to overfitting.

iv) Existence and uniqueness may not be guaranteed.

Example: MLE for a univariate Gaussian.

Setup: Given $D := \{x_1, x_2, \dots, x_n\}$, $x_i \in \mathbb{R}$



Assume $\left\{ \begin{array}{l} x_i \stackrel{i.i.d.}{\sim} p_\theta(x) = \mathcal{N}(x|\mu, \sigma^2) \quad , \quad \theta := \{\mu, \sigma^2\} \\ \Rightarrow x_i = \mu + \varepsilon \quad , \quad \varepsilon \sim \mathcal{N}(0, \sigma^2) \quad , \quad \sigma^2 > 0 \end{array} \right.$

generative model

Likelihood, $p(\mathcal{D}|\theta) = p(x_1, \dots, x_n | \mu, \sigma^2)$

$$\stackrel{i.i.d}{=} \prod_{i=1}^n p(x_i | \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right)$$

$$p(\mathcal{D}|\theta) := \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) := \mathcal{L}(\theta)$$

$$\theta_{MLE} := \arg \max_{\theta} p(\mathcal{D}|\theta) = \arg \min_{\theta} -\log p(\mathcal{D}|\theta)$$

$$-\log p(\mathcal{D}|\theta) = \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Take gradients and set them to zero to compute the critical points of $-\log p(\mathcal{D}|\theta) := \mathcal{L}(\mu, \sigma^2)$

$$\frac{\partial \mathcal{L}}{\partial \mu} = 0 \Rightarrow \frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - \mu) = 0 \Rightarrow \sum_{i=1}^n x_i - n\mu =$$

$$\Rightarrow \boxed{\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i}$$

Confirm that this is a minimum!

$$\frac{\partial^2 \mathcal{L}}{\partial \mu^2} = \frac{n}{\sigma^2} > 0 \quad , \quad \text{hence } \mu_{MLE} \text{ is a global and unique maximizer of the likelihood.}$$

Turn to optimization:

Setup: Given a model with parameters $\theta = (\theta_1, \dots, \theta_d)$ and a likelihood/loss $L(\theta)$ our goal is to estimate θ^* such that:

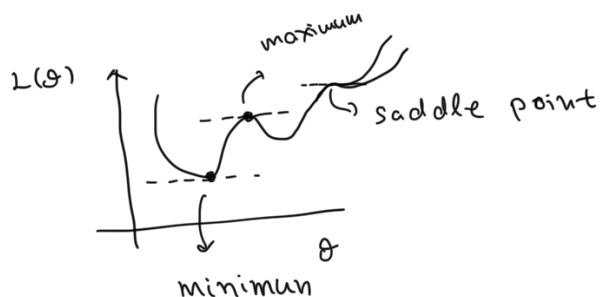
$$\boxed{\theta^* = \arg \min_{\theta} L(\theta)}, \quad L: \mathbb{R}^d \rightarrow \mathbb{R}$$

Gradient: $\nabla_{\theta} L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \\ \vdots \\ \frac{\partial L}{\partial \theta_d} \end{bmatrix}$
 $d \times 1$

We need to identify critical points:

$$\nabla_{\theta} L(\theta) = 0$$

This condition is met in 3 different scenarios:



Gradient descent:

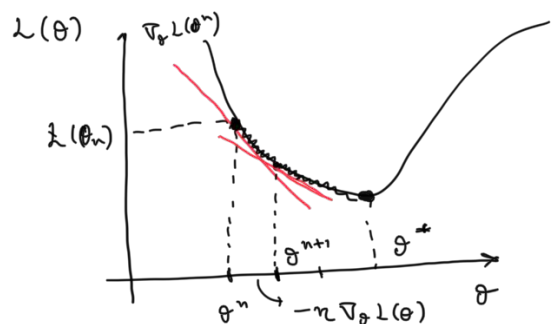
Starting from initial guess θ^n :

$$\rightarrow \theta^{n+1} = \theta^n - \eta \nabla_{\theta} L(\theta^n)$$

ODE
gradient
claw

$$\begin{cases} \frac{d\theta}{dt} = -\nabla_{\theta} L(\theta) \\ \dots \end{cases}$$

Forward-Euler
discretization



This is a first-order method

$$\theta(0) = \theta_0$$

... it relies on a linear approximation of $L(\theta)$ around θ .

η : step-size / learning rate

Hessian:

$$\nabla_{\theta}^2 L(\theta)_{d \times d} = \begin{bmatrix} \frac{\partial^2 L}{\partial \theta_1^2} & \frac{\partial^2 L}{\partial \theta_1 \partial \theta_2} & \dots & \frac{\partial^2 L}{\partial \theta_1 \partial \theta_d} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 L}{\partial \theta_d \partial \theta_1} & \dots & \dots & \frac{\partial^2 L}{\partial \theta_d^2} \end{bmatrix}$$

Taylor expansion of $L(\theta)$ around θ^n :

$$\hat{L}(\theta) \approx L(\theta^n) + g_n^T (\theta - \theta^n) + \frac{1}{2} (\theta - \theta^n)^T H_n (\theta - \theta^n) + \dots$$

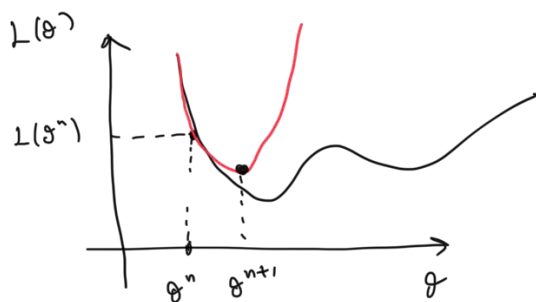
where $g_n := \nabla_{\theta} L(\theta^n)$, $H_n := \nabla_{\theta}^2 L(\theta^n)$.

Compute critical points:

$$\begin{cases} \nabla_{\theta} \hat{L}(\theta) = 0 \Rightarrow g_n^T + H_n \theta^{n+1} - H_n \theta^n = 0 \end{cases}$$

Second-order method.

$$\Rightarrow \boxed{\theta^{n+1} = \theta^n - H_n^{-1} g_n^T : \text{Newton's method}}$$



Pros: Utilizes the geometry of $L(\theta)$ better than gradient descent.

Cons: Computationally demanding for over-parametrized model. ($d \gg 1$)

e.g. in deep learning $d \sim \mathcal{O}(10^7)$