

Maximum Likelihood Estimation (MLE)

CHEN Si

Contents

1	Motivation	1
2	Maximum Likelihood Estimator	1
2.1	Interpretation	2
2.2	Properties	2
3	Conditional Log-Likelihood	3
3.1	Goal	3
3.2	Conditional Maximum Likelihood Estimator	3
3.3	Example: Linear Regression as Maximum Likelihood	3
	Appendix A Statistical Efficiency of MLE	5

1 Motivation

- Need some principle from which we can derive specific functions that are good estimators for different models

2 Maximum Likelihood Estimator

Problem Settings

- Consider a set of m examples $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ drawn **independently** from the **true but unknown** data-generating distribution $p_{\text{data}}(\mathbf{x})$
- Let $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ be a parametric family of probability distributions over the same space indexed by $\boldsymbol{\theta}$

Definition

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})\end{aligned}$$

- Disadvantage: Prone to numerical underflow

Logarithm of the Likelihood

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Rescale by m

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})]\end{aligned}$$

- This is an expectation with respect to the empirical distribution \hat{p}_{data} defined by the training data

2.1 Interpretation

- Minimize the cross-entropy

$$H(\hat{p}_{\text{data}}, p_{\text{model}}) = -\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

between the empirical distribution defined by the training set and the probability distribution defined by model

- Equivalently, minimize the dissimilarity between \hat{p}_{data} and \hat{p}_{model}

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})]$$

, whose minimum is 0

- Make the model distribution match the empirical distribution

2.2 Properties

It can be shown to be the best estimator **asymptotically**, as the number of examples $m \rightarrow \infty$, in terms of its rate of convergence.

1. Consistency, if under appropriate conditions:
 - The true distribution p_{data} must lie within the model family $p_{\text{model}}(\cdot; \boldsymbol{\theta})$
Otherwise, no estimator can recover p_{data}
 - The true distribution p_{data} must correspond to exactly one value of $\boldsymbol{\theta}$
Otherwise, maximum likelihood can recover the correct p_{data} but will not be able to determine which value of $\boldsymbol{\theta}$ was used by the data-generating process
2. Statistical Efficiency
 - The Cramér-Rao lower bound (Rao, 1945; Cramér, 1946) shows that no consistent estimator has a lower MSE than the maximum likelihood estimator (Appendix A)

For its consistency and efficiency, maximum likelihood is often considered the preferred estimator to use of machine learning.

When the number of examples is small enough to yield overfitting behavior, regularization strategies such as weight decay may be used to obtain a biased version of maximum likelihood that has less variance when training data is limited

3 Conditional Log-Likelihood

3.1 Goal

- Estimate a conditional probability $P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$
- In order to predict \mathbf{y} given \mathbf{x} (This forms the basis for most supervised learning)

3.2 Conditional Maximum Likelihood Estimator

Notations

- \mathbf{X} : All our inputs
- \mathbf{Y} : All our observed targets

Definition

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta})$$

If the examples are assumed to be **i.i.d.**, then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

3.3 Example: Linear Regression as Maximum Likelihood

Problem Settings

- The model produces a conditional distribution $p(y \mid \mathbf{x})$
- We define $p(y \mid \mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$
- The function $\hat{y}(\mathbf{x}; \mathbf{w})$ gives the prediction of the mean of the Gaussian
- We assume that the variance is fixed to some constant σ^2 chosen by the user

Mean Squared Error Since the examples are assumed to be **i.i.d.**, the conditional log-likelihood is given by

$$\sum_{i=1}^m \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \mathbf{w}) = -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}$$

Thus, maximizing the log-likelihood with respect to \mathbf{w} yields the same estimate of the parameters \mathbf{w} as minimizing the mean squared error

$$\text{MSE}_{\text{train}} = \frac{1}{m} \sum_{i=1}^m \|\hat{y}^{(i)} - y^{(i)}\|^2$$

Remark MSE is the cross-entropy between the empirical distribution defined by the training set and a Gaussian model defined by model

Appendix A Statistical Efficiency of MLE

There are other inductive principles besides the maximum likelihood estimator, many of which share the property of being consistent estimators. Consistent estimators can differ, however, in their statistical efficiency, meaning that one consistent estimator may obtain lower generalization error for a fixed number of samples n , or equivalently, may require fewer examples to obtain a fixed level of generalization error.

Statistical efficiency is typically studied in the parametric case (as in linear regression), where our goal is to estimate the value of a parameter (assuming it is possible to identify the true parameter), not the value of a function. A way to measure how close we are to the true parameter is by the expected mean squared error, computing the squared difference between the estimated and true parameter values, where the expectation is over training samples from the data-generating distribution. That parametric mean squared error decreases as n increases, and for n large, the Cramér-Rao lower bound (Rao, 1945; Cramér, 1946) shows that no consistent estimator has a lower MSE than the maximum likelihood estimator.