S&DS 242/542: Theory of Statistics Lecture 19: Parameter estimation in misspecified models "All models are wrong, but some are useful."

— George Box, 1976

Model misspecification

In our examples thus far, we typically started with the assumption $X_1, \ldots, X_n \stackrel{ID}{\sim} f(x \mid \theta)$. Here $f(x \mid \theta)$ plays two different roles:

- It describes the model we choose to fit to the data
- It describes the true distribution of the data

In reality these can be different, for various reasons:

- The true distribution of the data is unknown, and any model we fit is only an approximation to the truth
- The data is well-described by a complex model, but we choose to fit a simpler model that has more interpretable parameters
- We fit a model that is only intended for black-box prediction, rather than accurately describing the data-generating process

Q: How may we interpret the fitted parameter $\hat{\theta}$ and model $f(x \mid \hat{\theta})$ if $f(x \mid \theta)$ does not truly describe the distribution of the data?

Interpreting estimates in misspecified models

Suppose $(X_1, Y_1), \ldots, (X_n, Y_n) \in \mathbb{R}^2$ are covariate-response pairs, which are IID with some unknown distribution g(x, y).

We fit a linear regression model $Y \sim \beta_0 + \beta_1 X$ with parameters $\beta = (\beta_0, \beta_1)$ using least squares:

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname*{arg\,min}_{\beta_0, \beta_1} \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 X_i - Y_i)^2$$

Even if the true relation between X and Y is not linear, we may fit this model and interpret $(\hat{\beta}_0, \hat{\beta}_1)$ as representing the "best linear predictor" for Y based on X. For large *n*, we expect $(\hat{\beta}_0, \hat{\beta}_1)$ to be close to the minimizers of the population squared error

$$\arg\min_{\beta_0,\beta_1} \mathbb{E}_{(X,Y)\sim g}[(\beta_0+\beta_1X-Y)^2]$$

Interpreting estimates in misspecified models

Suppose $X_1, \ldots, X_n \stackrel{IID}{\sim} g(x)$, for an unknown distribution g(x).

We fit a model $f(x \mid \theta)$ with parameters $\theta \in \mathbb{R}^k$ using the method of moments. By definition, the estimated parameters $\hat{\theta}$ are those for which

$$\mathbb{E}_{\hat{\theta}}[X] = \frac{1}{n} \sum_{i=1}^{n} X_i, \quad \dots, \quad \mathbb{E}_{\hat{\theta}}[X^k] = \frac{1}{n} \sum_{i=1}^{n} X_i^k$$

This may be understood as estimating the distribution $f(x | \theta)$ within our model that matches g(x) in its first k moments.

For large *n*, we expect $\hat{\theta}$ to be close to the value of θ such that

$$\mathbb{E}_{\theta}[X] = \mathbb{E}_{X \sim g}[X], \quad \dots, \quad \mathbb{E}_{\theta}[X^k] = \mathbb{E}_{X \sim g}[X^k]$$

Interpreting estimates in misspecified models

Suppose $X_1, \ldots, X_n \stackrel{ID}{\sim} g(x)$, for an unknown distribution g(x). We fit a model $f(x \mid \theta)$ over parameters $\theta \in \Omega$ using maximum likelihood. What is the interpretation of $\hat{\theta}$?

Theorem

Let $D_{\text{KL}}(g(x) \parallel f(x \mid \theta))$ be the Kullback-Leibler (KL) divergence from $f(x \mid \theta)$ to g(x), and suppose that

$$\theta \mapsto D_{\mathrm{KL}}(g(x) \,\|\, f(x \mid \theta))$$

has a unique minimizer $\theta^* \in \Omega$. Then, under regularity conditions for $f(x \mid \theta)$, the MLE $\hat{\theta}$ converges to θ^* in probability as $n \to \infty$.

For large *n*, the MLE estimates the "KL projection" of the true distribution g(x) onto our model $\{f(x \mid \theta) : \theta \in \Omega\}$.

Kullback-Leibler divergence

Kullback-Leibler divergence

For two discrete distributions with PMFs f and g on a sample space \mathcal{X} , the **Kullback-Leibler (KL) divergence** from f to g is

$$D_{\mathrm{KL}}(g\|f) = \sum_{x \in \mathcal{X}} g(x) \log \frac{g(x)}{f(x)}$$

For two continuous distributions with PDFs f and g on the real line, the KL divergence from f to g is similarly

$$D_{ ext{KL}}(g\|f) = \int_{-\infty}^{\infty} g(x) \log rac{g(x)}{f(x)} \, dx$$

In both cases, this may be written as an expectation

$$D_{ ext{KL}}(g\|f) = \mathbb{E}_g\left[\log rac{g(X)}{f(X)}
ight] = \mathbb{E}_g[\log g(X)] - \mathbb{E}_g[\log f(X)]$$

where \mathbb{E}_g indicates expectation with respect to $X \sim g(x)$.

KL divergence between normal distributions

Let $f = \mathcal{N}(\mu_0, \sigma^2)$ and $g = \mathcal{N}(\mu_1, \sigma^2)$ with common variance σ^2 . To compute $D_{\text{KL}}(g \| f)$: $f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-M_0)^2}{2\sigma^2}} g(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-M_1)^2}{2\sigma^2}}$ => 1. g g(x) = - (x - m)2 + (x - m) $= \left(-\frac{1}{2\sigma^{1}} \times^{2} + \frac{\mu_{1}}{\sigma^{1}} \times - \frac{\mu_{1}^{2}}{2\sigma^{1}}\right) + \left(\frac{1}{2\sigma^{1}} \times^{2} - \frac{\mu_{2}}{\sigma^{1}} \times + \frac{\mu_{0}^{2}}{2\sigma^{1}}\right)$ $= \underbrace{\mu_{i}}_{m^{2}} \underbrace{\mu_{0}}_{m^{2}} \times - \underbrace{\mu_{i}}_{m^{2}} \underbrace{\mu_{0}}_{m^{2}} \times$ DKL (g 1 F) = Eg [M-M. X - M. H. 7

KL divergence between normal distributions

= /1.-10. IE (X) - Mi-10. = (m, (m, -m)) z So Dyn (glif) = (11-74)2

KL divergence between Bernoulli distributions

Let f = Bernoulli(p), g = Bernoulli(q). To compute $D_{\text{KL}}(g||f)$:



=> D_K (glf) = E [1., 2(x) /

= 1 log & + (1-1) log 17

Note that this expression is *not* symmetric in (p, q), so $D_{\text{KL}}(g||f) \neq D_{\text{KL}}(f||g)$.

KL divergence between Bernoulli distributions

For f = Bernoulli(p) and g = Bernoulli(q), we have

$$D_{ ext{KL}}(g\|f) = q\lograc{q}{p} + (1-q)\lograc{1-q}{1-p}$$

For p close to q, this may be approximated by Taylor expansion:

$$\begin{array}{l} \log p \approx \log q + (p^{-}q) \cdot \frac{1}{q} + \frac{(p^{-}q)^{2}}{2} \left(-\frac{1}{2^{2}}\right) \\ \Rightarrow \log q - \log p \approx -\frac{p^{-}q}{2} + \frac{(p^{-}q)^{2}}{2q^{2}} \\ \log (1-p) \approx \log (1-q)^{+} \left(p^{-}q\right) \cdot \left(-\frac{1}{1-q}\right) + \frac{(p^{-}q)^{2}}{2} \left(-\frac{1}{(1-q)^{2}}\right) \\ \Rightarrow \log (1-q)^{-1} \log (1-q)^{+} \left(p^{-}q\right) \cdot \left(-\frac{1}{1-q}\right) + \frac{(p^{-}q)^{2}}{2} \left(-\frac{1}{(1-q)^{2}}\right) \\ \Rightarrow \log (1-q)^{-1} \log (1-q) \approx \frac{p^{-}q}{1-q} \approx \frac{(p^{-}q)^{2}}{2((-q)^{2}} \end{array}$$

KL divergence between Bernoulli distributions

Dxc (g lf) = q (log q - log ((-1) (log ((-1) - log ((-1))) $\approx \frac{(p-1)^{2}}{2q} + \frac{(p-1)^{2}}{2((-\gamma))}$ $= \frac{(p-q)^{2}}{2q(l-q)} \approx \frac{(p-q)^{2}}{2p(l-p)}$ S. Er p≈q. Dru (gll E) ~ (p-2) ~ (p-2) L

KL divergence between Binomial distributions

Let f = Binomial(n, p) and g = Binomial(n, q). Then

$$\log \frac{g(x)}{f(x)} = \log \left(\binom{n}{x} q^x (1-q)^{n-x} \middle/ \binom{n}{x} p^x (1-p)^{n-x} \right)$$
$$= x \log \frac{q}{p} + (n-x) \log \frac{1-q}{1-p}$$

Applying $\mathbb{E}_{g}[X] = nq$,

$$\begin{split} D_{\mathrm{KL}}(g\|f) &= \mathbb{E}_g[X] \log \frac{q}{p} + (n - \mathbb{E}_g[X]) \log \frac{1 - q}{1 - p} \\ &= n \left(q \log \frac{q}{p} + (1 - q) \log \frac{1 - q}{1 - p} \right) \\ &= n \times D_{\mathrm{KL}}(\mathrm{Bernoulli}(q)\| \, \mathrm{Bernoulli}(p)) \end{split}$$

For p close to q, this is $D_{\mathrm{KL}}(g\|f) \approx n \times \frac{(p-q)^2}{2q(1-q)}$.

Properties of KL divergence

 $D_{\mathrm{KL}}(g \| f)$ is an information-theoretic measure of discrepancy between g and f. In general, $D_{\mathrm{KL}}(f \| g) \neq D_{\mathrm{KL}}(g \| f)$.

It satisfies the following basic properties:

- ▶ If f = g, then $D_{KL}(g||f) = 0$ because $\log \frac{g(x)}{f(x)} = 0$ for all x.
- D_{KL}(g || f) ≥ 0 for any f and g. This follows from applying Jensen's inequality to the convex function log x:

$$D_{\mu\nu}(g|F) = F_{g}\left[l_{g} \frac{g(\omega)}{f(\omega)}\right]$$

$$= E_{g}\left[-l_{0}g + \frac{f(x)}{g(x)}\right]$$

$$= -l_{0}g = E_{g}\left[\frac{f(x)}{g(x)}\right] = -l_{0}g + \int_{-\infty}^{\infty} \frac{f(x)}{g(x)} \frac{g(x)}{g(x)} \frac{dx}{dx}$$

$$= 0$$

KL divergence within a parametric model

Consider a parametric model $f(x \mid \theta)$ and two parameters $\theta_0, \theta_1 \in \mathbb{R}$. Then

$$D_{\mathrm{KL}}(f(x \mid \theta_0) \parallel f(x \mid \theta_1)) = \underbrace{\mathbb{E}_{\theta_0}[\log f(X \mid \theta_0)]}_{=L(\theta_0)} - \underbrace{\mathbb{E}_{\theta_0}[\log f(X \mid \theta_1)]}_{=L(\theta_1)}$$

Fixing θ_0 , let $L(\theta) = \mathbb{E}_{\theta_0}[\log f(X \mid \theta)]$ be the population log-likelihood function when the true parameter is θ_0 . If θ_0 and θ_1 are close to each other, then we may apply the Taylor expansion

$$L(heta_1)pprox L(heta_0)+(heta_1- heta_0)L'(heta_0)+rac{1}{2}(heta_1- heta_0)^2L''(heta_0)$$

Recall that

KL divergence within a parametric model

Thus, for parameters θ_0 and θ_1 close to each other,

$$D_{\mathrm{KL}}(f(x \mid heta_0) \parallel f(x \mid heta_1)) pprox rac{I(heta_0)}{2} (heta_1 - heta_0)^2$$

The KL divergence is approximately the squared difference between the parameter values, scaled by one-half times the Fisher information.

This gives another interpretation of the Fisher information: It relates the KL divergence — an information-theoretic measure of the discrepancy between $f(x \mid \theta_0)$ and $f(x \mid \theta_1)$ — to the squared difference of the model parameters. The MLE in misspecified models

Interpretation of the MLE

Suppose $X_1, \ldots, X_n \stackrel{ID}{\sim} g(x)$, for an unknown distribution g(x). We fit a model $f(x \mid \theta)$ over parameters $\theta \in \Omega$ using maximum likelihood.

Theorem Let $D_{KL}(g(x) || f(x | \theta))$ be the Kullback-Leibler (KL) divergence from $f(x | \theta)$ to g(x), and suppose that

 $\theta \mapsto D_{\mathrm{KL}}(g(x) \,\|\, f(x \mid \theta))$

has a unique minimizer $\theta^* \in \Omega$. Then, under regularity conditions for $f(x \mid \theta)$, the MLE $\hat{\theta}$ converges to θ^* in probability as $n \to \infty$.

Proof sketch

By definition, the MLE $\hat{\theta}$ maximizes

$$\frac{1}{n}\ell_n(\theta) = \frac{1}{n}\sum_{i=1}^n \log f(X_i \mid \theta)$$

By the Law of Large Numbers, as $n \to \infty$, for each $\theta \in \Omega$ this converges in probability to the expected log-likelihood function

$$L(\theta) = \mathbb{E}_g[\log f(X \mid \theta)]$$

Here, \mathbb{E}_g is now the expectation with respect to the true distribution $X \sim g(x)$.

In Lecture 14, we assumed that $g(x) = f(x \mid \theta^*)$ for some true parameter $\theta^* \in \Omega$. Then $L(\theta)$ is maximized at θ^* , and this explained consistency of the MLE.

Proof sketch

If g(x) does not belong to our model $\{f(x \mid \theta) : \theta \in \Omega\}$, let us write

$$\begin{split} L(\theta) &= \mathbb{E}_g[\log f(X \mid \theta)] \\ &= \mathbb{E}_g[\log g(X)] - \mathbb{E}_g\left[\log \frac{g(X)}{f(X \mid \theta)}\right] \\ &= \mathbb{E}_g[\log g(X)] - D_{\mathrm{KL}}(g(x) \parallel f(x \mid \theta)) \end{split}$$

The first term $\mathbb{E}_{g}[\log g(X)]$ does not depend on θ . Thus the maximizer of $L(\theta)$ is the minimizer of $D_{\mathrm{KL}}(g(x) || f(x | \theta))$, which is assumed in the theorem to be a unique value $\theta^* \in \Omega$. So we expect, as $n \to \infty$,

$$\hat{\theta} \to \theta^*$$

Example: MLE in a misspecified Poisson model

Suppose we fit the model Poisson(λ) to observations X_1, \ldots, X_n whose true distribution is g(x).

To compute
$$D_{KL}(glu) \parallel \mathcal{E}(x|\lambda))$$
 when
 $\mathcal{E}(x|\lambda) = \frac{e^{-\lambda}\lambda^{x}}{x!}$ is the PMF of Poisson(λ):
 $D_{KL}(glu) \parallel \mathcal{E}(x|\lambda)) = IE_{g}\left[log \frac{g(x)}{\mathcal{E}(x|\lambda)}\right]$
 $= IE_{g}\left[log g(x)\right] - IE_{g}\left[-\lambda + Xlog \lambda - log (X!)\right]$
 $= IE_{g}\left[log g(x)\right] + \lambda - IE_{g}[x] \cdot log \lambda + IE_{g}[log (x!)]$

Example: MLE in a misspecified Poisson model

Sot
$$O^{2} = \int_{X} D_{KL} (g(y) \parallel f(x|\lambda))$$

= $1 - \frac{E_{2}[X]}{\lambda}$

So the MLE $\hat{\lambda}$ satisfies

$$\hat{\lambda} \to \lambda^* = \mathbb{E}_g[X]$$

in probability as $n \to \infty$. This coincides with the direct conclusion using the Law of Large Numbers and the explicit form $\hat{\lambda} = \bar{X}$.

Example: Asymptotic variance of the MLE

Continuing this example, since $\hat{\lambda} = \bar{X}$, by the CLT we have

$$\sqrt{n}(\hat{\lambda} - \lambda^*) \rightarrow \mathcal{N}(0, \operatorname{Var}_g[X])$$

in distribution as $n \to \infty$. So $\operatorname{Var}[\hat{\lambda}] \approx \frac{\operatorname{Var}_g[X]}{n}$ for large n.

If we estimate this variance based on the Fisher information in our assumed $Poisson(\lambda)$ model, we would obtain an incorrect estimate:

$$I(\lambda) = \operatorname{Var}_{\lambda} \left[\frac{\partial}{\partial \lambda} \log \frac{e^{-\lambda} \lambda^{X}}{X!} \right] = \operatorname{Var}_{\lambda} \left[\frac{X}{\lambda} - 1 \right] = \frac{1}{\lambda}$$

Thus, since $\hat{\lambda} \approx \lambda^* = \mathbb{E}_g[X]$ for large *n*, the usual plug-in estimate for $Var[\hat{\lambda}]$ would be

$$\frac{1}{nI(\hat{\lambda})} = \frac{\hat{\lambda}}{n} \approx \frac{\lambda^*}{n} = \frac{\mathbb{E}_g[X]}{n} \neq \frac{\operatorname{Var}_g[X]}{n}$$

Uncertainty quantification in misspecified models

- It is often meaningful to fit a parametric model f(x | θ) that is only an approximation to the true distribution of data.
- If maximum likelihood is used to fit the model, then θ̂ may be understood as estimating a parameter θ* representing the KL projection of the data distribution onto the assumed model.

• However $Var[\hat{\theta}]$ is, in general, **not** given by $\frac{1}{nI(\theta^*)}$.

Thus the usual model-based estimate of Var[$\hat{\theta}$] given by $\frac{1}{nl(\hat{\theta})}$ may be incorrect in a misspecified model, and the confidence interval $\hat{\theta} \pm z^{(\alpha/2)} \sqrt{\frac{1}{nl(\hat{\theta})}}$ may have incorrect coverage for θ^* .