

Likelihood Ratio for Linear Models

April 9, 2019

Likelihood ratio for Linear Models

- Recall the linear model:

$$x = \mathbf{H}\boldsymbol{\theta} + w$$

- $w = [w[0] \dots w[N-1]]^T$ is a vector of i.i.d. samples with joint pdf $w \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- \mathbf{H} is the known model matrix.
- We want to test the simple hypotheses:

$$H_0 : \boldsymbol{\theta} = \mathbf{0}$$

$$H_1 : \boldsymbol{\theta} = \boldsymbol{\theta}_1, \text{ known}$$

- Here we assume we know $\boldsymbol{\theta}$, whereas previously in class we have attempted to estimate it.

Likelihood Ratio

- The likelihood ratio is:

$$\frac{p(x|H_1)}{p(x|H_0)} = \frac{\frac{1}{(2\pi)^{N/2}\sqrt{\det\sigma^2\mathbf{I}}} \exp[-1/(2\sigma^2)(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)^T(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)]}{\frac{1}{(2\pi)^{N/2}\sqrt{\det\sigma^2\mathbf{I}}} \exp[-1/(2\sigma^2)\mathbf{x}^T\mathbf{x}]} \underset{H_1}{\gtrless} \lambda$$

where $\lambda = \frac{\pi_0 L(1|0)}{\pi_1 L(0|1)}$ is the bayesian threshold for the likelihood ratio test.

Likelihood Ratio

- The likelihood ratio is:

$$\frac{p(x|H_1)}{p(x|H_0)} = \frac{\frac{1}{(2\pi)^{N/2}\sqrt{\det\sigma^2\mathbf{I}}} \exp[-1/(2\sigma^2)(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)^T(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)]}{\frac{1}{(2\pi)^{N/2}\sqrt{\det\sigma^2\mathbf{I}}} \exp[-1/(2\sigma^2)\mathbf{x}^T\mathbf{x}]} \stackrel{H_1}{\gtrless} \lambda$$

where $\lambda = \frac{\pi_0 L(1|0)}{\pi_1 L(0|1)}$ is the bayesian threshold for the likelihood ratio test.

- This simplifies, after some work to:

$$\mathbf{T}(\mathbf{x}) = \mathbf{x}^T \mathbf{H}\boldsymbol{\theta}_1 / \sigma^2 \stackrel{H_1}{\gtrless} \ln \lambda + \frac{\boldsymbol{\theta}_1^T \mathbf{H}^T \mathbf{H} \boldsymbol{\theta}_1}{2\sigma^2}$$

Likelihood Ratio

- The likelihood ratio is:

$$\frac{p(x|H_1)}{p(x|H_0)} = \frac{\frac{1}{(2\pi)^{N/2} \sqrt{\det \sigma^2 \mathbf{I}}} \exp[-1/(2\sigma^2)(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)^T(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}_1)]}{\frac{1}{(2\pi)^{N/2} \sqrt{\det \sigma^2 \mathbf{I}}} \exp[-1/(2\sigma^2)\mathbf{x}^T \mathbf{x}]} \stackrel{H_1}{\gtrless} \lambda$$

where $\lambda = \frac{\pi_0 L(1|0)}{\pi_1 L(0|1)}$ is the bayesian threshold for the likelihood ratio test.

- This simplifies, after some work to:

$$\mathbf{T}(\mathbf{x}) = \mathbf{x}^T \mathbf{H}\boldsymbol{\theta}_1 / \sigma^2 \stackrel{H_1}{\gtrless} \ln \lambda + \frac{\boldsymbol{\theta}_1^T \mathbf{H}^T \mathbf{H} \boldsymbol{\theta}_1}{2\sigma^2}$$

- $\mathbf{T}(\mathbf{x})$ is a normal random variable, so for the NP test we need to find the means and variances under each hypothesis to get P_D and P_F .

Mean and Variance under H_0

- For H_0 we have $x = H\mathbf{0} + w = w$:

$$\mathbb{E}(T(x)|H_0) = \mathbb{E}\left(\frac{w^T H \theta_1}{\sigma^2}\right) = \mathbf{0}$$

$$\begin{aligned} \text{var}(T(X)|H_0) &= \mathbb{E}\left[\left(\frac{w^T H \theta_1}{\sigma^2}\right)^T \frac{w^T H \theta_1}{\sigma^2}\right] = \mathbb{E}\left(\frac{\theta_1^T H^T}{\sigma^2} w w^T \frac{H \theta_1}{\sigma^2}\right) = \\ &= \frac{\theta_1^T H^T}{\sigma^2} \mathbb{E}(w w^T) \frac{H \theta_1}{\sigma^2} = \frac{\theta_1^T H^T}{\sigma^2} \sigma^2 I \frac{H \theta_1}{\sigma^2} = \frac{\theta_1^T H^T H \theta_1}{\sigma^2} \end{aligned}$$

- So $p(T(x)|H_0) \sim \mathcal{N}(\mathbf{0}, \frac{\theta_1^T H^T H \theta_1}{\sigma^2})$

Probability of False Alarm

- Let $\lambda' = \ln \lambda + \frac{\theta_1^T H^T H \theta_1}{2\sigma^2}$
- Then the probability of false alarm is:

$$\begin{aligned}
 P_F &= Pr(\mathbf{T}(\mathbf{x}) > \lambda' | H_0) = \\
 &Pr\left(\underbrace{\frac{(\mathbf{T}(\mathbf{x}) - 0)\sigma^2}{\theta_1^T H^T H \theta_1}}_{N(0,1)} > \frac{(\lambda' - 0)\sigma^2}{\theta_1^T H^T H \theta_1}\right) = \\
 &1 - \Phi\left(\frac{(\ln \lambda + \frac{\theta_1^T H^T H \theta_1}{2\sigma^2})\sigma^2}{\theta_1^T H^T H \theta_1}\right) = Q\left(\frac{(\ln \lambda + \frac{\theta_1^T H^T H \theta_1}{2\sigma^2})\sigma^2}{\theta_1^T H^T H \theta_1}\right)
 \end{aligned}$$

where Φ is the cumulative distribution function for the standard normal distribution.

Mean and Variance under H_1

- For H_1 use $x = H\theta_1 + w$:

$$\mathbb{E}(T(x)|H_1) = \mathbb{E}\left(\frac{(H\theta_1 + w)^T H\theta_1}{\sigma^2}\right) = \frac{\theta_1^T H^T H\theta_1}{\sigma^2}$$

- The variance is a bit more complicated:

$$\text{var}(T(x)|H_1) = \mathbb{E}\left(\left(\frac{x^T H\theta_1}{\sigma^2} - \mathbb{E}\left(\frac{x^T H\theta_1}{\sigma^2}\right)\right)^T \left(\frac{x^T H\theta_1}{\sigma^2} - \mathbb{E}\left(\frac{x^T H\theta_1}{\sigma^2}\right)\right)\right) =$$

$$\mathbb{E}\left(\left(x - \mathbb{E}(x)\right)^T \frac{H\theta_1}{\sigma^2}\right)^T \left(x - \mathbb{E}(x)\right)^T \frac{H\theta_1}{\sigma^2}\right) =$$

$$\text{cov}(x)\left(\frac{H\theta_1}{\sigma^2}\right)^T \left(\frac{H\theta_1}{\sigma^2}\right) = \frac{\theta_1^T H^T H\theta_1}{\sigma^2} = \text{var}(T(x)|H_0)$$

- So under H_1 , $T(x) \sim N\left(\frac{\theta_1^T H^T H\theta_1}{\sigma^2}, \frac{\theta_1^T H^T H\theta_1}{\sigma^2}\right)$, which has a different mean and the same variance as H_0

- The Probability of Detection is given by:

$$P_D = Pr(\mathbf{T}(\mathbf{x}) > \lambda' | H_1) = Pr\left(\underbrace{\frac{(\mathbf{T}(\mathbf{x}) - \frac{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}{\sigma^2})\sigma^2}{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}}_{N(0,1)} > \frac{(\lambda' - \frac{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}{\sigma^2})\sigma^2}{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}\right) =$$

$$1 - \Phi\left(\frac{(\ln \lambda - \frac{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}{2\sigma^2})\sigma^2}{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}\right) = Q\left(\frac{(\ln \lambda - \frac{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}{2\sigma^2})\sigma^2}{\theta_1^T \mathbf{H}^T \mathbf{H} \theta_1}\right)$$

- This is almost the same as P_F , with one subtraction instead of addition.
- Large λ leads to higher P_D and P_F
- $P_D = \Phi^{-1}(\Phi(P_F) - \sqrt{\frac{\theta_1^T \theta_1}{\sigma^2}})$ is another way to express P_D in this case - proof is in Kay CH. 3 and 4.

Example: Sum of Sinusoids, Kay ex. 4.9

- $H_0 : x[n] = w[n]$
- $H_1 : \mathbf{x}[n] = a \cos(2\pi f_0 n) + b \sin(2\pi f_0 n) + w[n]$
- $\theta_1 = \begin{bmatrix} a \\ b \end{bmatrix}$
- $\mathbf{H} = \begin{bmatrix} 1 & 0 \\ \cos(2\pi f_0) & \sin(2\pi f_0) \\ \vdots & \vdots \\ \cos(2\pi f_0(N-1)) & \sin(2\pi f_0(N-1)) \end{bmatrix}$
- For this example, let $N = 2$.

Example: Sum of Sinusoids, Kay ex. 4.9

- The test statistic is given by:

$$\frac{1}{\sigma^2} \mathbf{x}^T \mathbf{H} \boldsymbol{\theta}_1 = \frac{1}{\sigma^2} \begin{bmatrix} x[0] & x[1] \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \cos(2\pi f_0) & \sin(2\pi f_0) \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} =$$

$$\frac{1}{\sigma^2} (ax[0] + x[1](a \cos(2\pi f_0) + b \sin(2\pi f_0))) =$$

$$\frac{1}{\sigma^2} (a(x[0] + x[1] \cos(2\pi f_0)) + b(\sin(2\pi f_0))) =$$

$$a\hat{a} + b\hat{b}$$

- Here \hat{a} and \hat{b} are the estimators of a and b based on the fourier series formula!
- Under H_0 , \hat{a} and \hat{b} are very small, so $\mathbf{T}(x)$ is small.
- Under H_1 , $\mathbf{T}(x) \approx a^2 + b^2$. This is proportional to the signal power.

Comments

- What happens when θ_1 is unknown?
- We have to use the MLE estimator in something called the **Generalized Likelihood Ratio Test**.
- The hypotheses are slightly different:

$$H_0 : A\theta = b$$

$$H_1 : A\theta \neq b$$

- This hypothesis test whether or not θ lives in a particular subspace (e.g. line, plane, hyperplane)
- A , b are known and form a consistent (solveable) set of equations.
- Then you apply similar formulas, but the covariance matrix of T changes because $\hat{\theta}$ is a function of x and therefore a random variable.
- P_D and P_F become χ^2 distributions.

Comments on Testing Coefficients

- Suppose $H_0 = \theta_0 = [0, \theta_1, \theta_2, \dots, \theta_p]$
- And $H_1 = \theta_1 = [\theta_1, \theta_2, \theta_3, \dots, \theta_p]$ is almost exactly the same.
- The formulas are slightly more complicated since the mean of H_0 is no longer 0, but we can still apply the LRT. This is testing whether or not the data in the first column of \mathbf{H} improves the model or not, and is a simple version of where the p-values next to coefficients in statistics software comes from.
- This is often done in Linear Regression - I will talk more about this case when we get to the Generalized LRT.
- In machine learning this is a type of **Feature Selection** - there are many ways to do this for linear regression.