

STA 250: STATISTICS

Notes 20. Testing Goodness of Fit

Book chapters: 10.1-2

1 How good is a model?

So far in our course we have always accepted a given statistical model to describe our data. Can we ascertain whether a model is adequate for the data we observe? Surprisingly, Pearson's chi-square test provides quite a compelling way of doing it.

Let's start with a simpler problem. Suppose we model $Z = (Z_1, \dots, Z_n)$ as $Z_i \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$. Once we observe data $z = (z_1, \dots, z_n)$ can we ascertain if they really came from a $\text{Normal}(0, 1)$ distribution. We have done this in our labs many times. We simply make a histogram plot of our data and then overlay with the pdf of the $\text{Normal}(0, 1)$. If the histogram and the pdf look similar, we're reasonably happy that the $\text{Normal}(0, 1)$ model fits the observations.

In matching the histogram with the theoretical pdf, we essentially did the following. We split the range of our observations into bins, then counted number of observations in each bin and checked whether these counts are in agreement with the expected counts in the bins if they really came from the $\text{Normal}(0, 1)$ distribution. That is, we first turned our original data into a count data (with multinomial pmf) and then checked whether they agree with the bin probabilities under the $\text{Normal}(0, 1)$ pdf.

2 Formalization through bin count and Pearson's test

We can do this testing formally as follows. Fix k bins $(a_0 = -\infty, a_1], (a_1, a_2], \dots, (a_{k-2}, a_{k-1}], (a_{k-1}, a_k = \infty)$ covering the whole of $(-\infty, \infty)$, and define X_1 as the number of Z_i 's in the first bin, X_2 as the number of Z_i 's in the second bin and so on. Together we get k bin counts, and $X = (X_1, \dots, X_n)$ must have a $\text{Multinomial}(n, p)$ where p_l denotes the probability that Z_1 (or any other Z_i) would be in the l -th bin $(a_{l-1}, a_l]$.

If indeed $Z_i \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$ then we should have $p = p_0$ where $p_{0l} = \Phi(a_l) - \Phi(a_{l-1})$ for $l = 1, \dots, k$. Here $\Phi(z)$ denotes the $\text{Normal}(0, 1)$ cdf. So we can rephrase our question as testing $H_0 : p = p_0$ against $H_1 : p \neq p_0$. If we reject H_0 , then that would mean that we have evidence against Z_i 's being IID $\text{Normal}(0, 1)$.

Once the problem is set up this way, testing H_0 is straightforward by using Pearson's chi-square test for a point null. That is, we calculate $Q(x) = \sum_{l=1}^k (x_l - e_l)^2 / e_l$ where $e_l = np_{0l} = n(\Phi(a_l) - \Phi(a_{l-1}))$ and reject H_0 at level α if $Q(x) > F_{k-1}^{-1}(1 - \alpha)$, or report the p-value $1 - F_{k-1}(Q(x))$.

The only non-trivial thing here is how we choose the bins (how many and where). More bins would mean a finer comparison between the histogram shape and the $\text{Normal}(0, 1)$ pdf,

and hence better power. But by choosing too many bins, we may bring down the count in each to very small numbers, making the test unreliable. It is usually recommended that each bin should have expected count (under the null) of at least 5. So we can take k to be the integer just smaller than or equal to $n/5$.

We saw in the lab that the power of a Pearson's chi-square test against a point null is larger when the point null is $p_0 = (1/k, \dots, 1/k)$. So we can choose a_l 's to make the bins receive equal probabilities under the $\text{Normal}(0, 1)$ pdf. This is achieved by setting $a_1 = \Phi^{-1}(1/k)$, $a_2 = \Phi^{-1}(2/k)$, \dots , $a_{k-1} = \Phi^{-1}((k-1)/k)$. With this choice each $e_l = n/k$ so $Q(x) = \sum_{l=1}^k (x_l - n/k)^2 / (n/k) = \frac{k}{n} \sum_{l=1}^k (x_l - n/k)^2$.

Similarly, if we wanted to test if $Z_i \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$ for some fixed μ and σ^2 then we would do exactly as above, but choose the bins according to the inverse of the $\text{Normal}(\mu, \sigma^2)$ cdf, i.e., $a_1 = \mu + \sigma\Phi^{-1}(1/k)$, $a_2 = \mu + \sigma\Phi^{-1}(2/k)$, \dots , $a_{k-1} = \mu + \sigma\Phi^{-1}((k-1)/k)$. Once we get the bin counts x_1, \dots, x_k we have a similar construction of the test statistic $Q(x) = \frac{k}{n} \sum_{l=1}^k (x_l - n/k)^2$. As before, $Q(X)$ is approximately $\chi^2(k-1)$ under the null. Therefore a size α test rejects H_0 if $Q(x) > F_{k-1}^{-1}(1-\alpha)$ and the p-value based on such tests is $1 - F_{k-1}(Q(x))$.

3 Testing goodness-of-fit of the normal model

However, our real interest is in checking if $Z_i \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$ for some μ and σ^2 that are not necessarily known to us. To tackle this we can follow the same recipe above but with μ and σ estimated by $\hat{\mu} = \bar{z}$ and $\hat{\sigma} = s_z$. That is, we use bins $a_1 = \bar{z} + s_z\Phi^{-1}(1/k)$, $a_2 = \bar{z} + s_z\Phi^{-1}(2/k)$, etc to get our bin counts x_1, \dots, x_k . We again construct $Q(x) = \frac{k}{n} \sum_{l=1}^k (x_l - n/k)^2$.

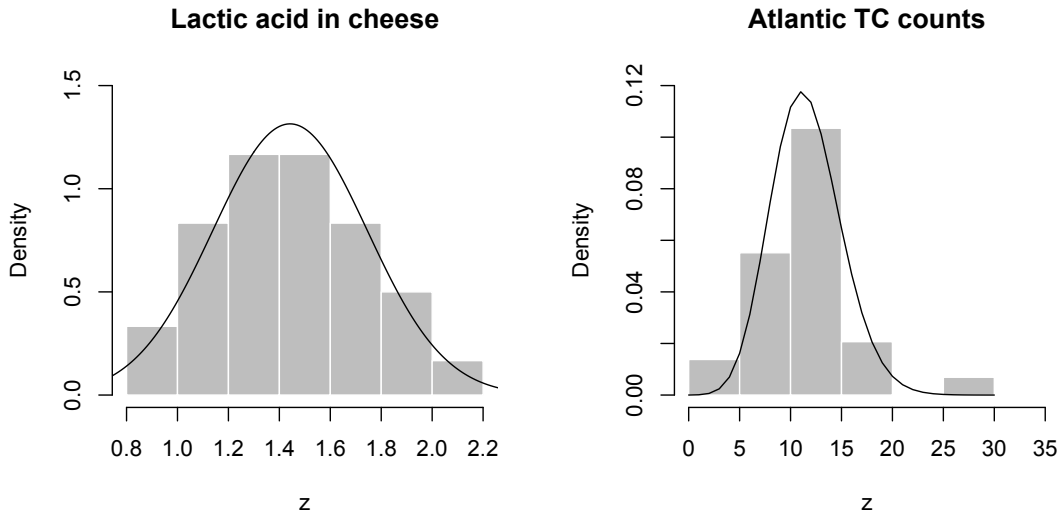
What should be the distribution of this $Q(X)$ under the null? The number of free parameters under the null is $r = 2$ because we had both μ and σ^2 unknown. Should we have $Q(X) \sim \chi^2(k-1-2) = \chi^2(k-3)$? This approximation, however, does not work so well. But it turns out that under the null $P(Q(X) > c)$ lies between $1 - F_{k-3}(c)$ and $1 - F_{k-1}(c)$. Consequently the p-value based on "reject H_0 if $Q(x) > c$ " type of tests lies between $1 - F_{k-3}(Q(x))$ and $1 - F_{k-1}(Q(x))$. So we have a range of p-values.

Example (Lactic acid concentration). 30 samples from a cheese slab are measured to have the following lactic acid concentrations:

0.86, 1.53, 1.57, 1.81, 0.99, 1.09, 1.29, 1.78, 1.29, 1.58,
 1.68, 1.90, 1.06, 1.30, 1.52, 1.74, 1.16, 1.49, 1.63, 1.99,
 1.15, 1.33, 1.44, 2.01, 1.31, 1.46, 1.72, 1.25, 1.08, 1.25

We want to test whether these observations are IID draws from a $\text{Normal}(\mu, \sigma^2)$ distribution, with no assumption made on μ and σ^2 . The figure below (left) shows the data histogram overlaid with the $\text{Normal}(\hat{\mu}, \hat{\sigma}^2)$ pdf where $\hat{\mu} = \bar{z} = 1.44$ and $\hat{\sigma} = s_z = 0.30$. There seems a fairly good agreement between the two.

Because we have $n = 30$, we can use $k = n/5 = 6$ bins of the form $(a_{l-1}, a_l]$ with $a_l = \bar{z} + s_z\Phi^{-1}(l/k)$. The bins are found to be $(-\infty, 1.15]$, $(1.15, 1.31]$, $(1.31, 1.44]$, $(1.44, 1.57]$, $(1.57, 1.74]$, $(1.74, \infty)$ with bin counts 5, 8, 2, 5, 4, 6. We calculate $Q(x) = 4$. Therefore the p-value ranges between $1 - F_3(4) = 0.26$ and $1 - F_5(4) = 0.55$. So we would fail to accept the null hypotheses at nominal levels of 5%, 10% or 1%.



4 Testing goodness-of-fit of a general parametric model

The same concepts apply for a general model $Z_i \stackrel{\text{iid}}{\sim} g(z_i|\theta)$, $\theta \in \Theta$ for data $Z = (Z_1, \dots, Z_n)$. We first find an estimate $\hat{\theta}(z)$ of θ , usually the MLE, and then form $k \geq n/5$ bins $(a_{l-1}, a_l]$, $l = 1, \dots, k$ with $a_l = G^{-1}(l/k|\theta = \hat{\theta}(x))$, where $G(z_i|\theta)$ denotes the CDF of $g(z_i|\theta)$. We get bin counts x_1, \dots, x_k and construct $Q(x) = \frac{k}{n} \sum_{l=1}^k (x_l - n/k)^2$ and find a p-value range between $1 - F_{k-1-r}(Q(x))$ and $1 - F_{k-1}(Q(x))$ where r is the dimension of θ (i.e., number of free parameters we needed to estimate).

Example (Annual TC counts in the north Atlantic). The annual tropical cyclone counts in the north Atlantic between 1980 and 2009 are as below (30 years).

11 11 5 4 12 11 6 7 12 11
 14 8 6 8 7 19 13 7 14 12
 14 15 12 16 14 27 10 14 16 11

Are these data from a $\text{Poisson}(\mu)$ distribution for some $\mu > 0$? To test this, we first estimate μ by $\hat{\mu}_{\text{MLE}}(z) = 11.57$. The figure above (right) shows the data histogram overlaid with the $\text{Poisson}(11.57)$ pmf. A total of $n/5 = 6$ bins are found by using the inverse CDF of this Poisson pmf. These equal $[0, 8]$, $(8, 10]$, $(10, 11]$, $(11, 13]$, $(13, 15]$, $(15, \infty)$ with bin counts 9 1 5 5 6 4. For these counts, we get $Q(x) = 6.8$. Consequently the p-value ranges between $1 - F_4(6.8) = 0.147$ to $1 - F_5(6.8) = 0.236$. So again there is no strong evidence against the null hypothesis.