

①

Nonparametric Goodness-of-Fit Tests

We've seen in Chapter 10 how to test whether or not proportions (p_1, \dots, p_k) were equal to prescribed ones $(p_{1,0}, \dots, p_{k,0}) = p_0$. Testing for $H_0: p = p_0$ amounts to test if the unknown distribution P of data equals Multinomial (m, p_0) , or equivalently $H_0: P = \text{Multinomial}(m, p_0)$.

In other words, this is a goodness-of-fit test with

- 1) Discrete distributions (since we deal with counts/bins here)
- 2) A parametric model (namely, the multinomial distributions)

In these notes, we explore the idea above for

- 1) Continuous distribution
- 2) Nonparametric models.

(2)

Notation: let $X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P$ be real random variables.

Since we work in a nonparametric setting, we don't have parameters/formulas to describe distributions. Hence, we must work in full generality with its cumulative distribution function

$$F_x(t) = F(t) = P(X \leq t), \quad t \in \mathbb{R}.$$

Recall that $F_x = F_y \iff X$ and Y have the same distribution, so cdf's are good candidates for describing/testing equality of distributions.

We have $F(t) = P(X \leq t) = E[\mathbb{1}_{\{X \leq t\}}]$, An estimator of

F is the empirical distribution function) Plug in / Method of moments

$$F_m(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{X_i \leq t\}}$$

Indeed, we see that for any fixed $t \in \mathbb{R}$, the strong law of large numbers yields

$$F_m(t) \xrightarrow[m \rightarrow \infty]{\text{a.s.}} F(t).$$

Actually, the convergence is much stronger: it's uniform.

(3)

Thm (Glivenko - Cantelli)

$$\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow \infty]{a.s.} 0$$

Proof: let $\varepsilon > 0$. Then, fix $k > \frac{1}{\varepsilon}$ and consider "knot" points x_0, \dots, x_k such that

$$-\infty = x_0 < x_1 \leq x_2 \leq \dots \leq x_{k-1} < x_k = +\infty$$

that define a partition of \mathbb{R} into k disjoint intervals such that

$$F(x_j^-) \leq \frac{j}{k} \leq F(x_j) \quad , \quad 1 \leq j \leq k-1$$

where $F(t^-) = \mathbb{P}(X < t) = F(t) - \mathbb{P}(X = t)$
 $= \lim_{s \uparrow t} F(s)$,

Then, by construction, if $x_{j-1} < x_j$,

$$F(x_j^-) - F(x_{j-1}) \leq \frac{j}{k} - \frac{(j-1)}{k} = \frac{1}{k} < \varepsilon .$$

(4)

By by the strong LLN, we have the pointwise convergences

$$F_n(x_j) \xrightarrow[n \rightarrow \infty]{a.s.} F(x_j) \text{ and } F_n(x_j^-) \xrightarrow[n \rightarrow \infty]{a.s.} F(x_j^-),$$

which rewrites as

$$|F_n(x_j) - F(x_j)| \xrightarrow[n \rightarrow \infty]{a.s.} 0 \text{ and } |F_n(x_j^-) - F(x_j^-)| \xrightarrow[n \rightarrow \infty]{a.s.} 0.$$

Hence, looking at the maximum over all j in the finite set $\{1, \dots, k-1\}$,

$$\Delta_n = \max_{1 \leq j \leq k-1} \left\{ |F_n(x_j) - F(x_j)|, |F_n(x_j^-) - F(x_j^-)| \right\} \xrightarrow[n \rightarrow \infty]{} 0.$$

Now, notice that any $t \in \mathbb{R}$ falls into some interval $[x_{j-1}, x_j)$.

$$x_{j-1} \leq t < x_j.$$

We then have:

$$F_n(t) - F(t) \leq F_n(x_j^-) - F(x_{j-1}) \leq F_n(x_j^-) - F(x_j^-) + \varepsilon$$

and

$$F_n(t) - F(t) \geq F_n(x_{j-1}) - F(x_j) \geq F_n(x_{j-1}) - F(x_{j-1}) - \varepsilon,$$

(5)

and thus for any $t \in \mathbb{R}$,

$$|F_m(t) - F(t)| \leq \Delta_m + \varepsilon \xrightarrow[m \rightarrow \infty]{a.s.} \varepsilon$$

As this holds for arbitrary $t \in \mathbb{R}$, it follows that

$$\forall \varepsilon > 0, \limsup_{m \rightarrow \infty} \sup_{t \in \mathbb{R}} |F_m(t) - F(t)| \leq \varepsilon \quad a.s.$$

All that remains to be done is to send ε to 0. For this, write

$$\text{the event } A_\varepsilon = \left\{ \limsup_{m \rightarrow \infty} \sup_{t \in \mathbb{R}} |F_m(t) - F(t)| \leq \varepsilon \right\}.$$

Clearly, $\varepsilon' < \varepsilon \Rightarrow A_{\varepsilon'} \subset A_\varepsilon$, and $A_0 = \bigcap_{\varepsilon > 0} A_\varepsilon$, so that

$$P(A_0) = P\left(\bigcap_{\varepsilon > 0} A_\varepsilon\right) = \inf_{\varepsilon > 0} P(A_\varepsilon)$$

But from what's above, $P(A_\varepsilon) = 1$ for all $\varepsilon > 0$. Hence, $P(A_0) = 1$, or equivalently

$$P\left(\lim_{m \rightarrow \infty} \sup_{t \in \mathbb{R}} |F_m(t) - F(t)| = 0\right) = 1$$

⑥

Rk: Although Glivenko - Cantelli is a very strong result, it does not help us derive a testing procedure

Kolmogorov - Smirnov Test for Goodness of Fit

We introduce a (pseudo)-distance between the empirical distribution $P_n = \sum_{i=1}^n \delta_{X_i}$ and P , that simply is the L^∞ distance between their respective cdf's:

$$D_{KS}(P_n, P) = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)|$$

If F is given, this is a statistic that can actually be computed straightforwardly.

To see this, write $(X_{(1)}, \dots, X_{(n)})$ for the order statistic of (X_1, \dots, X_n) . Namely it's just the sample points sorted by increasing order

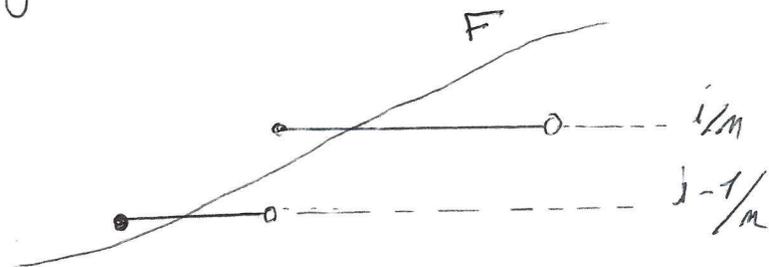
$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

(7)

Prop: $D_{KS}(P_n, P) = \max_{1 \leq i \leq n} \max \left\{ \left| F(X_{(i)}) - \frac{i}{n} \right|, \left| F(X_{(i)}) - \frac{i-1}{n} \right| \right\}$

Proof: F is increasing, and F_n is piecewise constant and taking values

$\left\{ \frac{i}{n} \right\}_{0 \leq i \leq n}$



We now describe the distribution of $D_{KS}(P_n, P)$.

Thm: Write $D_n = \sqrt{n} D_{KS}(P_n, P) = \sup_{t \in \mathbb{R}} \sqrt{n} |F_n(t) - F(t)|$.

Assume that P has a density

(i) If $U_n = \sup_{t \in [0,1]} \sqrt{n} |F_{U,n}(t) - t|$, where $F_{U,n}$ is the empirical distribution function associated to a iid n -sample of the uniform distribution on $[0, 1]$, then D_n has the same distribution as U_n . (rich tables are known)

(ii) (Kolmogorov, 1933) D_n converges in distribution towards the distribution with cdf

$$H(t) = 1 + 2 \sum_{k=1}^{\infty} (-1)^k e^{-2k^2 t^2}$$

(iii) If P_0 is a distribution with density, and that $P \neq P_0$, then

$$D_n^{\circ} = \sqrt{n} D_{KS}(P_n, P_0) \xrightarrow[n \rightarrow \infty]{P} \infty$$

⑧ To prove the theorem, we'll need some results on the quantile function

$$F^{-1}(u) = \text{Inf} \{t \mid F(t) \geq u\}$$

Also called
generalized inverse.

Lemma F and F^{-1} satisfy

- 1) $\forall t \in \mathbb{R}, F^{-1}(F(t)) \leq t$
- 2) $\forall u \in (0, 1), F(F^{-1}(u)) \geq u$
- 3) $F^{-1}(u) \leq t \iff u \leq F(t)$
- 4) If $U \sim \text{Uniform}([0, 1])$, $F^{-1}(U)$ has the same distribution as X .

Proof: let us write $\mathcal{T}_u = \{t \mid F(t) \geq u\}$, so that $F^{-1}(u) = \text{Inf} \mathcal{T}_u$.

1) $\forall t \in \mathbb{R}, t \in \mathcal{T}_{F(t)}$, so $F^{-1}(F(t)) \leq t$.

2) By definition, there exists a decreasing sequence $(t_n)_n$ of elements of \mathcal{T}_u such that $t_n \xrightarrow{n \rightarrow \infty} F^{-1}(u)$. But since $F(t_n) \geq u$ and that F is right continuous, $F(F^{-1}(u)) \geq u$.

3) If $u \leq F(t)$, then $t \in \mathcal{T}_u$ so that $F^{-1}(u) \leq t$.

conversely
If $F^{-1}(u) \leq t$, then since F is non-decreasing, $F(F^{-1}(u)) \leq F(t)$.
Using 2), we get $F(t) \geq u$.

9

4) For $t \in \mathbb{R}$, from 3) above,

$$\begin{aligned} P(F^{-1}(U) \leq t) &= P(U \leq F(t)) \\ &= F(t), \end{aligned}$$

and since the cdf characterises fully the distribution of X , we get the result. \square

Proof of the theorem

(i) If $U_i \sim \text{Uniform}([0, 1])$, $F^{-1}(U_i) \stackrel{\mathcal{D}}{=} X_i$. As a consequence,

$$\begin{aligned} F_{U_{(m)}}(F(t)) &= \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{U_i \leq F(t)} \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{F^{-1}(U_i) \leq t} \\ &\stackrel{\mathcal{D}}{=} F_m(t). \end{aligned}$$

Taking suprema over $t \in \mathbb{R}$ we get that $D_m \stackrel{\mathcal{D}}{=} \sqrt{m} \sup_{t \in \mathbb{R}} |F_{U_{(m)}}(F(t)) - F(t)|$.

Since F is continuous (because P has a density), $\{F(t)\}_{t \in \mathbb{R}} = (0, 1)$, so that we get (by change of variable $u = F(t)$), that $U_m \stackrel{\mathcal{D}}{=} D_m$.

(ii) Not proved here. See Billingsley for instance.

(10)

(iii) If $P_0 \neq P$, there exists $\varepsilon > 0$ and $t \in \mathbb{R}$ such that $|F(t) - F_0(t)| \geq \varepsilon$.
Without loss of generality, assume that $F(t) - F_0(t) \geq \varepsilon$.

Then,

$$D_m^0 \geq \sqrt{m} |F_m(t) - F_0(t)|$$

Triangle inequality $\left\{ \begin{array}{l} \downarrow \\ \leftarrow \end{array} \right.$

$$\geq \sqrt{m} |F(t) - F_0(t)| - \sqrt{m} |F(t) - F_m(t)|$$

$$\geq \sqrt{m} \varepsilon - \sqrt{m} |F(t) - F_m(t)|.$$

But from the LLN, $|F_m(t) - F(t)| \xrightarrow[m \rightarrow \infty]{P} 0$, so

$$P(|F_m(t) - F(t)| \geq \frac{\varepsilon}{2}) \xrightarrow[m \rightarrow \infty]{} 0,$$

which yields $P(D_m^0 \geq \frac{\sqrt{m} \varepsilon}{2}) \xrightarrow[m \rightarrow \infty]{} 1$,

and hence $D_m^0 \xrightarrow[m \rightarrow \infty]{P} \infty \left| \left(\forall \eta > 0, P(D_m^0 \geq \eta) \xrightarrow[m \rightarrow \infty]{} 1 \right) \right.$

(11)

Application of these results to building the Goodness-of-Fit K.S. Test

let P_0 be a fixed (known) distribution, with cdf F_0 that is continuous.

Consider the testing problem

$$H_0: P = P_0 \quad \text{vs} \quad H_1: P \neq P_0$$

Test Statistic: Considering the previous results, we take

$$\begin{aligned} D_m^0 &= \sqrt{m} D_{KS}(P_m, P_0) \\ &= \sqrt{m} \text{Max}_{1 \leq i \leq m} \text{Max} \left\{ \left| F_0(X_{(i)}) - \frac{i}{m} \right|, \left| F_0(X_{(i)}) - \frac{i-1}{m} \right| \right\} \end{aligned}$$

Rejection Region: Since $D_m^0 \xrightarrow{P} \infty$ under H_1 , we take

$$R = \{ D_m^0 \geq \lambda \} \quad \text{for some } \lambda \in (0, 1)$$

Calibration of λ :

- If m is small ($m \leq 100$), we use the tables of the distribution of U_n
- If m is large ($m > 100$), we use those of the asymptotic distribution with cdf H .

12

Rk: The K.S. test has both nonasymptotic and asymptotic formulations

• The K.S. test cannot be used if P_0 does not have a density!
(i.e. F_0 continuous)

• Astonishingly, we often note that the K.S. test is more powerful than the Goodness-of-fit χ^2 test, especially if n is small

• As for χ^2 test, one can adapt K.S. to test if P belongs to some parametric family $\{P_\theta, \theta \in \Theta\}$: uniform distributions, Gaussians, Exponentials ...

In these cases one needs to show that by taking an estimator $\hat{\theta}$ of θ , the distribution of $\sqrt{n} D_{KS}(P_n, P_{\hat{\theta}})$ does not depend on P .

We'll see later on an example of these for Gaussians (Lilliefors)

(13)

Kolmogorov-Smirnov Test for Homogeneity

We consider the two-sample problem that consists in testing if two samples have the same distribution.

Namely, we have two independent sequences

$$\begin{cases} X_1, \dots, X_m \sim_{\text{iid}} P_X \\ Y_1, \dots, Y_m \sim_{\text{iid}} P_Y \end{cases}$$

where P_X and P_Y are unknown.

We let $F_m(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{X_i \leq t}$ and $G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{Y_j \leq t}$ denote

the empirical distribution functions of the X_i 's and Y_j 's respectively.

Write

$$D_{m,m} = \sqrt{\frac{m}{m+m}} \sup_{t \in \mathbb{R}} |F_m(t) - G_m(t)|$$

Thm Assume that P_X and P_Y have a density.

(i) The distribution of $D_{m,m}$ does not depend on P_X and P_Y when $P_X = P_Y$. It only depends on m and m .

(ii) If $P_X \neq P_Y$, then $D_{m,m} \xrightarrow[m, m \rightarrow \infty]{P} \infty$

(14)

Application of this result to building the K-S. test for homogeneity

We want to test

$$H_0: P_X = P_Y \quad \text{vs} \quad H_1: P_X \neq P_Y$$

Test statistic: $D_{m,m} = \sqrt{\frac{m \cdot m}{m+m}} \sup_{t \in \mathbb{R}} |F_m(t) - G_m(t)|$

Rejection Region: $\mathcal{R} = \{D_{m,m} \geq \Delta\}$ for some $\Delta > 0$

Calibration of Δ : Using the tables of the K-S distribution with parameters (m, m)

Rk: The idea of comparing cdf and empirical cdf is pretty general.

Instead of taking the sup-norm between these two functions, one can also consider:

• L^2 distance: $m \int_{\mathbb{R}} (F_m(t) - F(t))^2 f(t) dt$ (Cramer-von Mises test, see Homework 6)

• weighted L^2 distance: $m \int_{\mathbb{R}} \frac{(F_m(t) - F(t))^2}{F(t)(1-F(t))} f(t) dt$ (Anderson-Darling test)

⊕ two-sample versions of these.

Lilliefors Test of Normality

Back to goodness-of-fit tests, we sometimes want to test whether data comes from a family of distributions (as opposed to a single P_0 of interest).

For instance, one may wonder if P is a Gaussian / Exponential / Beta / Gamma / ... distribution.

Namely, given $X_1, \dots, X_m \stackrel{iid}{\sim} P$ and a parametric model $\{P_\theta \mid \theta \in \Theta\}$, you'd like to test

$$H_0: P \in \{P_\theta \mid \theta \in \Theta\} \quad \text{vs} \quad H_1: P \notin \{P_\theta \mid \theta \in \Theta\}.$$

Among these, probably the most important family is the Gaussian model $\{N(\mu, \sigma^2) \mid \mu \in \mathbb{R}, \sigma > 0\}$.

The Lilliefors test of Normality roughly boils down (⚠ a distribution with random parameters do not exist!) to apply a goodness-of-fit K.S. test to a distribution \hat{P}_0 is Gaussian with mean $\bar{X}_m = \frac{1}{m} \sum_{i=1}^m X_i$ and variance $S_m^2 = \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X}_m)^2$.

Hypotheses: $H_0: P$ is Gaussian $H_1: P$ is not Gaussian

test statistic: let F_0 denote the cdf of the distribution $N(0, 1)$.

We let $\hat{F}_0(t) = F_0\left(\frac{t - \bar{X}_n}{S_n}\right)$, and define the test statistic by

$$L_n = \sqrt{n} \sup_{t \in \mathbb{R}} |F_n(t) - \hat{F}_0(t)|$$

$$= \sqrt{n} \max_{1 \leq i \leq n} \max \left\{ \left| F_0\left(\frac{X_{(i)} - \bar{X}_n}{S_n}\right) - \frac{i}{n} \right|, \left| F_0\left(\frac{X_{(i)} - \bar{X}_n}{S_n}\right) - \frac{i-1}{n} \right| \right\}$$

Rejection Region: $R = \{L_n \geq \lambda\}$ for some $\lambda > 0$

Calibrations: If H_0 holds (meaning $P = N(\mu, \sigma^2)$ for some (μ, σ)), the distribution of $\frac{X_i - \bar{X}_n}{S_n}$ does not depend on (μ, σ) , so the distribution of L_n does not depend on P under H_0 .

⚠ The distribution of L_n under H_0 is not the K.S. distribution! So λ is to be read in another table: that of the Lilliefors distribution.