# Handout 11a - Towards distribution-freeness

## 1  Empirical processes indexed by functions

So far, we have have seen that one can perform goodness-of-fit for univariate distributions considering the empirical process

$$v_n(x; F) = \sqrt{n}[F_n(x) - F(x)] \tag{1}$$

which can be seen as a process indexed by intervals of the type $(-\infty, x]$. Whereas, to test multivariate distributions, we rely on the (multivariate) empirical process

$$v_n(\boldsymbol{x}; F) = \sqrt{n}[F_n(\boldsymbol{x}) - F(\boldsymbol{x})], \quad \boldsymbol{x} = (x_1, \ldots, x_p) \tag{2}$$

which can be seen as a process indexed by the sets $(-\infty, x_1] \times (-\infty, x_2] \times \ldots (-\infty, x_p]$. All of this can be further generalized by considering the empirical processes indexed by functions. For our purposes, we will consider functions $\phi \in \mathcal{L}(F) \subset L^2(F)$. Where $L^2(F)$ is the set of square integrable functions w.r.t. $F$, i.e.,

$$\phi: \quad \underbrace{\int \phi_{(x)}^2 \, dF(x)}_{f(x)\,dx} = <\phi, \phi>_F = E_F\left[\phi^2\right] < +\infty$$

whereas $\mathcal{L}(F)$ is a subset of $L^2(F)$ which collects functions which not only are square integrable w.r.t $F$ but are also orthogonal to 1 or, in other words, have mean zero under $F$, i.e., Finally,

$$\int \phi(x) \cdot 1 \, dF(x) = <\phi, 1>_F = E_F\left[\phi\right] = 0$$

the empirical process indexed by $\phi \in \mathcal{L}(F)$ is defined as

$$V_n(\phi; F) = \int \phi(x) \, dV_n(x; F) \qquad \text{STOCHASTIC INTEGRAL}$$

$$= \int \phi(x) \, d\sqrt{n}\left[F_n(x) - F(x)\right]$$

$$= \sqrt{n}\left[\int \phi(x) \, dF_n(x) - \int \phi(x) \, dF(x)\right] \longrightarrow = 0$$

$$= \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \phi(x_i) \qquad \frac{1}{n}\sum_{i=1}^{n} \phi(x_i)$$

Let's rewrite the multivariate empirical process in equation 2 in this form:

$$V_n(\theta; F) = \sqrt{n}\left[F_n(x) - F(x)\right]$$

$$= \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\left[\underbrace{\mathbb{1}_{\{x_i \leq x\}} - F(x)}_{=\phi_x(x_i)}\right] = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_x(x_i) = V_n(\phi_x; F)$$

- **What are its mean and covariance function?**

$$E_F\left[\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_x(x_i)\right] \underset{\text{i.i.d}}{=} \sqrt{n}\, E_F\left[\phi_x(x_1)\right] = \sqrt{n}\,<\phi_x, 1>_F = 0$$

$$E_F\left[\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_x(x_i)\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\phi_{x'}(x_i)\right] = \frac{1}{n}\left\{\sum_{i=1}^{n}E_F\left[\phi_x(x_i)\phi_{x'}(x_i)\right] + \right.$$

$$\left. +2\sum_{i'<i}\underbrace{E_F\left[\phi_x(x_i)\phi_{x'}(x_{i'})\right]}_{=0}\right\}$$

$$= <\phi_x, \phi_{x'}>_F$$

$$= F(x \wedge x') - F(x)F(x') < +\infty$$

$$\Rightarrow \phi_x \in \mathcal{L}(F)$$

- **What is the limiting distribution of $v_n(\phi_x; F)$?**

$$F - \text{Brownian Bridge}$$

which is a Gaussian process... and thus it is fully characterized by its mean and covariance $<\phi_x, 1>_F$ and $<\phi_x, \phi_{x'}>_F$.

## 2    Constructing distribution free tests

All the consideration above hold for any arbitrary distribution function. To perform distribution-free goodness-of-fit we will consider a _reference distribution Q_, and the respective empirical process:

$$V_n(\psi_x; Q) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\psi_x(x_i) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\left[\underbrace{\mathbb{1}_{\{x_i \leq x\}} - Q(x)}_{\psi_x(x_i)}\right]$$

with mean $<\psi_x, 1>_Q = 0$

covariance $<\psi_x, \psi_{x'}>_Q < +\infty$        $\psi_x \in \mathcal{L}(Q)$

2

$$\widetilde{\phi}_x(x_i)$$

*But how do we use it?* Instead of testing $F$ directly by taking functionals of $v_n(\phi_x; F)$ (see Handout 11), we will construct another empirical process, namely $v_n(l\psi_x; F)$ (to be defined in a minute) which, under $F$, has the same limiting distribution of $v_n(\psi_x; Q)$, under $Q$. Our test statistics will correspond to functionals of $v_n(l\psi_x; F)$ and, under F, will have the same distribution as functionals of $v_n(\psi_x; Q)$, under $Q$. This can be done by exploiting the so-called *Khmaladze-2 (K-2) transform* and which consists of the following steps:

**Step 1** - Map $\psi_x$ into $L^2(F)$ via the isometry

$$\mathcal{L}(Q)$$

$$l(t) = \sqrt{\frac{q(t)}{f(t)}},$$

where $q$ and $f$ are the densities of $Q$ and $F$, respectively. Obtain $l\psi_x \in L^2(F)$. To see that:

$$\langle l\psi_x, l\psi_x \rangle_F = \int l^2(t)\psi_x^2(t)\underbrace{dF(t)}_{f(t)dt} = \int \frac{q(t)}{f(t)}\psi_x^2(t)f(t)dt$$

$$= \int \psi_x^2(t)q(t)dt = \int \psi_x^2(t)dQ(t) = \langle \psi_x, \psi_x \rangle_Q < +\infty$$

$$l\psi_x \in L^2(F)$$

Do the $l\psi_x \in \mathcal{R}(F.)$?

$$\langle l\psi_x, 1 \rangle_F = \int l(t)\psi_x(t)dF(t)$$

$$= \int \sqrt{\frac{q(t)}{f(t)}}\psi_x(t)f(t)dt$$

$$= \int \psi_x(t)\sqrt{f(t)q(t)}\,dt \neq 0$$

$$\Rightarrow l\psi_x \notin \mathcal{R}(F)$$

$\in L^2(F)$

**Step 2 -** Map $l\psi_x$ into $\mathcal{L}(F)$ by means of the the unitary operator

$$K = I - \frac{1-l}{1-\langle l,1\rangle_F}\langle 1-l,\cdot\rangle_F,$$  (3)

and obtain $Kl\psi_x \in \mathcal{L}(F)$. To see that:

it preserves the inner product

$$\langle Kx, Ky\rangle = \langle x, y\rangle$$

$Kl\psi_x = K\Big(l(t)\,\psi_x(t)\Big)$

$= l(t)\,\psi_x(t) - \dfrac{1-l(t)}{1-\int l(t)\,dF(t)}\int(1-l(t))\,l(t)\,\psi_x(t)\,dF(t)$

$= l(t)\,\psi_x(t) - \dfrac{1-l(t)}{1-\int l(t)\,dF(t)}\Big[\int l(t)\,\psi_x(t)\,dF(t) - \underbrace{\int l^2(t)\,\psi_x(t)\,dF}\Big]$

$\underbrace{\phantom{\int l^2(t)\,\psi_x(t)\,dF}} = \langle l^2\psi, 1\rangle_F$

$= \langle \psi, 1\rangle_Q = 0$

$= l(t)\,\psi_x(t) - \dfrac{1-l(t)}{1-\int l(t)\,dF(t)}\int l(t)\,\psi_x(t)\,dt$

is it in $L^2(F)$? yes

$\langle Kl\psi_x, Kl\psi_x\rangle_F = \langle l\psi_x, l\psi_x\rangle_F = \langle \psi_x, \psi_x\rangle_Q < +\infty$

is $Kl\psi$ in $\mathcal{L}(F)$?

$\langle Kl\psi_x, 1\rangle_F = \int l\psi_x\,dF - \dfrac{1-\int l\,dF}{1-\int l\,dF}\int l\psi_x\,dF = 0$

4

It follows that $v_n(Kl\psi_x; F)$ and $v_n(\psi_x; Q)$, converge to a Gaussian process $v(\psi_x; Q)$ with mean and covariance:

MEAN $\quad \langle kl\psi_x, 1 \rangle_F = \langle \psi_x, 1 \rangle_Q = 0$

COVARIANCE $\quad \langle kl\psi_x, kl\psi_{x'} \rangle_F = \langle l\psi_x, l\psi_{x'} \rangle_F = \langle \psi_x, \psi_{x'} \rangle_Q < +\infty$

Hence we construct an entire family of test statistics for testing $H_0 : X \sim F$ versus $H_1 : X \nsim F$, i.e.,

- **Kolmogorov's statistics:** $\sup |v_n(Kl\psi_x; F)| \xrightarrow{d} \sup |v(\psi_x; Q)|$.

- **Cramer von Mises statistics:** $\int |v_n(Kl\psi_x; F)|^2 dQ(x) \xrightarrow{d} \int |v(\psi_x; Q)|^2 dQ(x)$.

- **Anderson-Darling statistics:** $\int \left| \frac{v_n(Kl\psi_x; F)}{\sqrt{Q(x)(1-Q(x))}} \right|^2 dQ(x) \xrightarrow{d} \int \left| \frac{v(\psi_x; Q)}{\sqrt{Q(x)(1-Q(x))}} \right|^2 dQ(x)$.

where the convergence is intended as $n \to \infty$, under $H_0$. So, for sufficiently large $n$, a valid testing procedure consists of simulating the distribution of the functionals of $v_n(\psi_x; Q)$ under $Q$ and using it to assign significance to the values of the functionals of $v_n(Kl\psi_x; F)$ observed on the data. For instance, if we decide to use Kolmogorov's statistics, we simulate the distribution of $D_Q = \sup_{\psi_x} |v_n(\psi_x; Q)|$ under $Q$. This will give us our null distribution. We then evaluate evaluate $D_{obs} = \sup_{\psi_x} |v_n(Kl\psi_x; F)|$ on the data observed. Our p-value is:

$$P(D_Q \geq D_{obs})$$