

INFINITE OCCUPANCY

A. D. BARBOUR

Uni. Zürich

CLASSICAL OCCUPANCY

n balls thrown independently and uniformly into
 m boxes: $N_j := \#$ balls in box j , $1 \leq j \leq m$

$$\sum_{j=1}^m N_j = n.$$

Questions: what can be said about

$K_{n0} := \#$ empty boxes;

$K_{nr} := \#$ boxes w/ exactly r balls;

Sequentially: how many balls thrown until the
first time that $n-r$ boxes are occupied, T_r :

Coupon collector problem.

Unequal box probabilities:

box i chosen w. pr. p_i , $1 \leq i \leq m$: $\sum_{i=1}^m p_i = 1$

Sampling theory.

$N \sim MN(n; p_1, \dots, p_m)$; various non-linear

functions of N are of interest.

INFINITE OCCUPANCY

Bahadur 1960

Karlin 1967

$p_1 \geq p_2 \geq \dots$, $\sum_{j \geq 1} p_j = 1$, all positive.

$K := K_n := \#$ occupied boxes

$$= \sum_{j \geq 1} \mathbb{I}[N_j \geq 1] =: \sum_{j \geq 1} I_j;$$

$K_{n,r}$, $r \geq 1$, as before.

Sequentially it is clear that $K_n \rightarrow \infty$ a.s.

Karlin showed: $K_n / \mathbb{E}K_n \rightarrow 1$ a.s.

CLT for K_n ?

Ex: $p_j = q^{j-1}(1-q)$, $j \geq 1$, some $0 < q < 1$:

Geometric probabilities. Then, since $N_j \sim \text{Bi}(n, p_j)$,

we have

$$\mathbb{E}I_j = P[N_j \geq 1] \sim 1 - e^{-np_j}; \quad \text{var } I_j \sim e^{-np_j}(1 - e^{-np_j});$$

and, since $(N_j, j \geq 1)$ are negatively associated,

$$\text{var } K_n \leq \sum_{j \geq 1} \text{var } I_j \sim \sum_{j \geq 1} e^{-np_j}(1 - e^{-np_j}).$$

Define j_n to satisfy $nq^{j_n-1}(1-q) \geq 1 > nq^{j_n}(1-q)$.

Then

$$\sum_{j \geq 1} e^{-np_j} (1 - e^{-np_j})$$

$$\leq np_{j_n+1} \sum_{s \geq 0} q^s + \sum_{s \geq 0} e^{-np_{j_n}} q^{-s}$$

$$\leq \sum_{s \geq 0} q^s + \sum_{s \geq 0} e^{-e^{-s}} < \infty.$$

So $\text{var } K_n$ remains bounded, and the CLT does not hold.

REGULAR VARIATION

Karlin defined a measure ν by

$$\nu(x, \infty) := \# \{j : p_j \geq x\}.$$

The sequence $(p_j, j \geq 1)$ is called regularly varying with exponent β , $0 < \beta < 1$, if

$$\nu(x, \infty) \sim l(1/x) x^{-\beta} \text{ as } x \downarrow 0,$$

for $l(\cdot)$ slowly varying.

[e.g. $p_j \sim j^{-1/\beta}$]. 3.

Karlin's Theorems: if regularly varying β , then

(i) (Sequentially) $K_n \sim T(1-\beta) n^\beta l(n)$ a.s.

(ii) $\mathbb{E}K_n \asymp \text{var} K_n \asymp n^\beta l(n)$

(iii) $(K_n - \mathbb{E}K_n) / \sqrt{\text{var} K_n} \rightarrow_d \mathcal{N}(0,1)$

(iv) $(K_{nr} - \mathbb{E}K_{nr}) / \sqrt{n^\beta l(n)}$ are asymptotically jointly ($r \geq 1$) normal as $n \rightarrow \infty$.

Dutko (1989) showed that

$(K_n - \mathbb{E}K_n) / \sqrt{\text{var} K_n} \rightarrow_d \mathcal{N}(0,1)$ iff $\text{var} K_n \rightarrow \infty$.

WARNING.

For $V_r(n) := \text{Var} K_{n,r}$, can have

• $\limsup_{n \rightarrow \infty} V_r(n) = \infty$ for all r ;

• $V_r(n) \rightarrow \infty$ for all $r \leq r_0$;

• $\liminf_{n \rightarrow \infty} V_r(n) = 0$ for all $r > r_0$;

} Simultaneous

and

• joint multivariate normal convergence ONLY IF R.V.

POISSONIZATION

If instead $\boxed{\text{Po}(n)}$ balls are thrown, and $L_j := \#$ balls in box j , then L_1, L_2, \dots are independent, with $L_j \sim \text{Po}(np_j)$; so the limit theory should not be too surprising.

Cheap Poisson argument.

Choose any sequence $j_n \rightarrow \infty$. Then

$$P_n := \sum_{j \geq j_n} p_j \rightarrow 0.$$

Thus $N_n^+ = \sum_{j \geq j_n} N_j$ has mean nP_n , and

$$d_{TV}(\mathcal{L}(N_n^+), \text{Po}(nP_n)) \leq P_n \rightarrow 0 \quad (n \rightarrow \infty).$$

Le Cam - Michel observation:

$$d_{TV}(\mathcal{L}(N_j, j \geq j_n), \mathcal{L}(L_j, j \geq j_n)) \leq P_n \rightarrow 0$$

also, for $(L_j, j \geq j_n)$ independent $\text{Po}(np_j)$ -distributed r.v.'s

Now choose $j_n \rightarrow \infty$ so slowly that

$$p_{j_n-1} \geq \frac{4 \log n}{n} > p_{j_n}.$$

$$\begin{aligned} \text{Then } P\left[\bigcup_{j < j_n} \{N_j = 0\}\right] &\leq (j_n - 1) \left(1 - \frac{4 \log n}{n}\right)^n \\ &\leq \frac{n}{4 \log n} \cdot n^{-4} \leq n^{-3} \rightarrow 0 \end{aligned}$$

as $n \rightarrow \infty$. Hence, defining

$$\tilde{K}_n := j_n - 1 + \sum_{j \geq j_n} \mathbb{I}[L_j \geq 1],$$

we have

$$d_{TV}(\mathcal{L}(K_n), \mathcal{L}(\tilde{K}_n)) \leq n^{-3} + P_n,$$

and K_n satisfies the CLT because \tilde{K}_n does,

if

$$\text{Var } \tilde{K}_n = \sum_{j \geq 1} e^{-np_j} (1 - e^{-np_j}) \rightarrow 0.$$

error rate $n^{-3} + P_n + 1/\sqrt{\text{Var } \tilde{K}_n}$.

LOCAL LIMIT THEOREMS

Hwang and Janson (Ann. Prob. 2008):

local limit approximation for K_n with error
of order $O(1/\text{var} K_n)$ if $\text{var} K_n \rightarrow \infty$.

E.g. if regularly varying β , then

$$\frac{1}{\text{var} K_n} \asymp \frac{1}{n^{\beta} \ell(n)}, \text{ whereas } P_n \asymp \left(\frac{\log n}{n}\right)^{1-\beta}.$$

So the cheap argument would give their result
(in R.V. case) ONLY for $\beta < 1/2$.

ROLLIN'S THEOREM for translated Poisson approximation.

Let W be an integer valued random variable, X

Some random element. Define

$$\mu_X := \mathbb{E}(W|X); \quad \mu = \mathbb{E}W = \mathbb{E}\mu_X;$$

$$\sigma_X^2 := \text{var}(W|X); \quad \rho^2 = \mathbb{E}\sigma_X^2;$$

$$\tau^2 := \text{var}\mu_X; \quad \nu^2 = \text{var}\sigma_X^2.$$

Suppose that

$$U := (\mu_X - \mu) / \tau$$

satisfies

$$|\mathbb{E}\{f'(U) - Uf(U)\}| \leq \varepsilon \|f''\| \quad (\otimes)$$

for nice f . Then

$$d_{\text{loc}}(\mathcal{L}(W), \text{TP}(\mu, \sigma^2))$$

$$\leq \mathbb{E} d_{\text{loc}}(\mathcal{L}(W|X), \text{TP}(\mu_X, \sigma_X^2))$$

$$+ O\left(\frac{1}{\rho^2} \left(1 + \frac{\nu^2}{\rho^2} + \frac{\varepsilon \tau^3}{\sigma^2}\right)\right),$$

where $\sigma^2 := \tau^2 + \rho^2 = \text{var}W$.

So if, for instance, $\mathcal{L}(W|X)$ is the law of a sum of independent indicators, and if $\rho^2 \asymp \sigma^2$, $V^2 = O(\rho^2)$ and $\varepsilon\tau^3 = O(\sigma^2)$, then a local limit approximation to $\mathcal{L}(W)$ of ideal order $O(1/\sigma^2)$ is obtained.

————— " —————

How do we use this here?

- j_n as before, so that $P[N_j \geq 1, \text{all } j < j_n] \geq 1 - n^{-3}$.
- j_0 fixed, with $P_0 := \sum_{j \geq j_0} p_j$ not too big.

Realize $(N_j, j \geq j_0) \sim MN(n; \underbrace{(p_j, j \geq j_0)}_{\text{defective dist}^n})$

by taking $(M_j, j \geq j_0) \sim MN(n; \underbrace{([p_j/P_0], j \geq j_0)}_{\text{probability dist}^n})$

and then 'thinning' independently, with retention probability P_0 .

Then $N_j \sim \text{Bi}(M_j, P_0)$ are conditionally independent given $M := (M_j, j \geq j_0)$, and hence so are the indicators $(\mathbb{I}[N_j \geq 1], j \geq j_n)$.

$$M_M = \sum_{j \geq j_n} \{1 - (1 - P_0)^{M_j}\} =: \sum_{j \geq j_n} q(M_j)$$

$$\sigma_M^2 = \sum_{j \geq j_n} q(M_j)(1 - q(M_j))$$

$$\rho^2 = \mathbb{E} \sigma_M^2 = \sum_{j \geq j_n} \mathbb{E} \{ (1 - P_0)^{M_j} - (1 - P_0)^{2M_j} \}$$

$M_j \sim \text{Bi}(n, p_j/P_0)$ gives explicit formulae:

$$\rho^2 \sim \sum_{j \geq j_n} e^{-np_j} (1 - e^{-np_j(1-P_0)})$$

$$\approx \sigma^2 \sim \sum_{j \geq j_n} e^{-np_j} (1 - e^{-np_j})$$

Then $\nu^2 = \text{var}(\sigma_M^2) = O(\sigma^2)$ also
(similar computation).

It thus remains to show that, for $U := (\mu_n - \mu) / \tau$,
we have

$$| \mathbb{E} \{ f'(U) - U f''(U) \} | \leq \varepsilon \| f'' \| ,$$

with $\varepsilon \tau^3 = o(\sigma^2)$, and we're done.

- extensions.