

Non-Gaussian local limit properties for pattern statistics in rational stochastic models

Massimiliano Goldwurm¹, Jianyi Lin², Marco Vignati¹

¹Dipartimento di Matematica, Università degli Studi di Milano, Italy

²Department of Mathematics, Khalifa University, Abu Dhabi, United Arab Emirates

AofA - CIRM - Marseille, 24 - 28 June 2019

[GLV18] Proc. DCFS 2018, 20 Int. Conf. Descriptive Complexity of Formal Systems, LNCS 10952.

[GLV19] Proc. DLT 2019, 23rd Int. Conf. Developments of Language Theory, LNCS 11647.

Pattern Statistics

finite alphabet	A
pattern	$a \in A,$ $w \in A^+, w = m$ $R \subseteq A^*,$ finite, regular, ...
random text	$x \in A^+, x = n$ stochastic model (Bernoullian, Markovian,..., rational)

$$O_n = \#\{\text{occurrences of pattern in } x\} \quad (\text{positions})$$

$$O_n \in \{0, 1, \dots, n\}$$

- Goals:
- asymptotic properties of $\{O_n\}$
 - $E(O_n), \text{var}(O_n), \dots$;
 - limit distributions, $\Pr\{O_n \leq z\} \sim \dots$
 - local limit laws, $\Pr\{O_n = k\} \sim \dots$

[Guibas-Odlyzko 78, 81], [Regnier-Szpankowski 96, 98],

[Nicodeme-Salvy-Flajolet 02], [Flajolet-Szpankowski-Vallée 06], ...

Rational models

[BCGL03, BCGL06, DGL04, GL06]

Rational formal series over $\{a, b\}$

$$r : \{a, b\}^* \rightarrow \mathbb{R}_+, \quad r = \sum_{w \in \{a, b\}^*} r(w)w$$

defined by a f.s. automaton with weights in \mathbb{R}_+

$$\Rightarrow L = \{x \in \{a, b\}^* \mid r(x) > 0\} \text{ is regular}$$

Symbol statistics ($\forall n \in \mathbb{N}$)

$$Y_n = |x|_a \quad \text{independent r.v.}$$

$$x \in \{a, b\}^n \text{ with } \Pr(x) = \frac{r(x)}{\sum_{w \in \{a, b\}^n} r(w)}$$

$$\Rightarrow Y_n \in \{0, 1, \dots, n\}, \quad Y_n = Y_n(r)$$

Special cases:

- $r = \chi_L \Rightarrow x \in L \cap \{a, b\}^n$ under uniform dist.
- Markovian source (π, P) with set of states A

$$\pi' P^n e = \sum_{w \in \{a, b\}^n} r(w) = 1, \quad \forall n \in \mathbb{N}$$

Regular pattern + Markovian source \subsetneq Rational models

[Nicodeme-Salvy-Flajolet 02, BCGL03]

1) For any:

- ▶ $R \subseteq A^*$: regular language (pattern)
- ▶ (π, P) : Markovian source over A

$\Rightarrow \exists r \in \mathbb{R}_+^{\text{rat}} \langle \langle a, b \rangle \rangle$ s.t.

$O_n(\pi, P, R)$ and $Y_n(r)$ have the **same distribution**

i.e. $\forall k = 0, \dots, n, \Pr(O_n = k) = \Pr(Y_n = k)$.

2) The opposite of 1) does not hold:

$$\sum_{n=0}^{+\infty} \sum_{k=0}^n \Pr(O_n = k) x^k y^n \text{ is a } \mathbf{rational} \text{ function } \quad \forall (\pi, P, R)$$

while for the rational model

$$\sum_{n=0}^{+\infty} \sum_{k=0}^n \Pr(Y_n = k) x^k y^n \text{ is } \mathbf{not rational} \text{ in general.}$$

Formal notions

Linear representation for $r \in \mathbb{R}_+^{\text{rat}} \langle\langle a, b \rangle\rangle$

(ξ, A, B, η) :

- $\xi, \eta \in \mathbb{R}_+^m$, weights of initial and final states
 - $A, B \in \mathbb{R}_+^{m \times m}$ ($\neq 0$), weights of a - and b -transitions
- A, B generate a morphism $\mu : \{a, b\}^* \rightarrow \mathbb{R}_+^{m \times m}$

$$\mu(a) = A, \mu(b) = B$$

$$\mu(a_1 a_2 \cdots a_n) = \mu(a_1) \mu(a_2) \cdots \mu(a_n)$$

Rational series $r \in \mathbb{R}_+^{\text{rat}} \langle\langle a, b \rangle\rangle$ defined by (ξ, A, B, η) :
for every $x = a_1 a_2 \cdots a_n \in \{a, b\}^*$

$$r(x) = \xi' \mu(x) \eta = \xi' \mu(a_1) \mu(a_2) \cdots \mu(a_n) \eta$$

Probability Measure over $\{a, b\}^n : \forall x \in \{a, b\}^n$

$$\Pr(x) = \frac{r(x)}{\sum_{w \in \{a, b\}^n} r(w)} = \frac{\xi' \mu(w) \eta}{\xi'(A + B)^n \eta}$$

special case: If $r = \chi_L$ where $L \subseteq \{a, b\}^*$ is **regular**, then

$$\Pr(x) = \begin{cases} \frac{1}{\#(L \cap \{a, b\}^n)} & \text{if } x \in L \\ 0 & \text{otherwise} \end{cases}$$

Symbol statistics

$$Y_n = |x|_a, \quad \text{where } \Pr(x) = \frac{r(x)}{\sum_{w \in \{a, b\}^n} r(w)}$$

$$Y_n \in \{0, 1, \dots, n\} \text{ is a r.v.,} \quad Y_n = Y_n(r)$$

Distribution of Y_n

$$\begin{aligned} \Pr(Y_n = k) &= \frac{\sum_{x \in \{a,b\}^n, |x|_a = k} r(x)}{\sum_{w \in \{a,b\}^n} r(w)} = \\ &= \frac{[z^k] \xi'(Az + B)^{n\eta}}{\xi'(A + B)^{n\eta}}, \quad \text{for } k = 0, \dots, n. \end{aligned}$$

Characteristic function of Y_n

$$\psi_n(t) = \frac{\xi'(Ae^{it} + B)^{n\eta}}{\xi'(A + B)^{n\eta}} = \frac{h_n(it)}{h_n(0)}$$

\implies asymptotic behaviour of $\{Y_n\}_n$

Case analysis according to $M = A + B$ [BCGL03, DGL04, BCGL06]

{	primitive	Gaussian	
	dominant bicomponent	Gaussian	
	equipotent bicomponent comm.	$\begin{cases} E_1 \neq E_2 & \text{uniform} \\ E_1 = E_2, V_1 \neq V_2 & \text{normal Mixture} \\ E_1 = E_2, V_1 = V_2 & \text{Gaussian} \end{cases}$	

Primitive rational models

[BCGL03, BCGL06]

(ξ, A, B, η) with primitive $M = A + B$, $A \neq 0 \neq B$

Main parameters: λ, β, γ

$\lambda > 0$ main eigenvalue of M

$v, u \in \mathbb{R}_+^m$ left - right normalized eigenvectors (λ)
 $v'M = \lambda v'$, $Mu = \lambda u$, $v'u = 1$, $v, u > 0$

$\beta = \frac{v'Au}{\lambda}$ mean constant, $0 < \beta < 1$

$\gamma > 0$ variance constant,
 $\gamma = \beta - \beta^2 + \frac{v'ADAu}{\lambda^2}$, $D \in \mathbb{R}^{m \times m} \rightsquigarrow$ other eigen.

Results:

$$E(Y_n) = \beta n + c + o(1), \quad c \in \mathbb{R}$$

$$\text{var}(Y_n) = \gamma n + O(1)$$

$$\frac{Y_n - \beta n}{\sqrt{\gamma n}} \rightarrow \mathcal{N}(0, 1) \quad \text{in dist.}$$

Local limit law (in the primitive case)

[BCGL06, GLV19]

Aperiodicity Condition :

G labelled graph with a -label weights A_{ij}
 b -label weights B_{ij}

$$d = \text{GCD}\{|C_1|_a - |C_2|_a : i \overset{C_1}{\sim} i, i \overset{C_2}{\sim} i, |C_1| = |C_2|\}$$

(A, B) **aperiodic** if $d = 1$

Theorem If M primitive, $A \neq 0 \neq B$, (A, B) aperiodic, then

$$\left| \sqrt{n} \Pr(Y_n = k) - \frac{e^{-\frac{(k-\beta n)^2}{2\gamma n}}}{\sqrt{2\pi\gamma}} \right| = O(n^{-1/2})$$

uniformly for every $k \in \{0, 1, \dots, n\}$.

Bicomponent models (with communication)

[dGL04]

Model formed by two communicating irreducible components
 (ξ, A, B, η) of size $m_1 + m_2$ such that

$$\xi' = (\xi'_1, \xi'_2), \quad A = \left[\begin{array}{c|c} A_1 & A_0 \\ \hline 0 & A_2 \end{array} \right], \quad B = \left[\begin{array}{c|c} B_1 & B_0 \\ \hline 0 & B_2 \end{array} \right], \quad \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix}$$

$(\xi_1, A_1, B_1, \eta_1)$ of size m_1 , $(\xi_2, A_2, B_2, \eta_2)$ of size m_2

$$M = \left[\begin{array}{c|c} A_1 + B_1 & A_0 + B_0 \\ \hline 0 & A_2 + B_2 \end{array} \right]$$

Hypotheses:

- 1) $A_1 + B_1$ and $A_2 + B_2$ irreducible
with λ_1, λ_2 main eigenvalues
- 2) $A_0 + B_0 \neq 0$ and $\xi_1 \neq 0 \neq \eta_2$
communication cond. $1 \sim > 2$

\Rightarrow various cases depending on: $\lambda_1 \stackrel{?}{=} \lambda_2, \beta_1 \stackrel{?}{=} \beta_2, \gamma_1 \stackrel{?}{=} \gamma_2$

Dominant bicomponent models: $\lambda_1 \neq \lambda_2$

Hp: $\lambda_1 > \lambda_2$, M_1 primitive, $A_1 \neq 0 \neq B_1$ ($\Rightarrow 0 < \beta_1 < 1$, $0 < \gamma_1$)

$$E(Y_n) \sim \beta_1 n, \quad \text{var}(Y_n) \sim \gamma_1 n, \quad \frac{Y_n - \beta_1 n}{\sqrt{\gamma_1 n}} \xrightarrow{d} \mathcal{N}(0, 1)$$

[dGL04]

Theorem Under the same hps, if (A_1, B_1) aperiodic then

$$\left| \sqrt{n} \Pr(Y_n = k) - \frac{e^{-\frac{(k - \beta_1 n)^2}{2\gamma_1 n}}}{\sqrt{2\pi\gamma_1}} \right| = O(n^{-1/2})$$

uniformly for every $k \in \{0, 1, \dots, n\}$.

[GLV19]

Equipotent bicomponent models with different β 'sHp: $\lambda_1 = \lambda_2$, $\beta_1 \neq \beta_2$, M_1, M_2 primitive, $A_j \neq 0 \neq B_j$ ($j = 1, 2$)

$$E(Y_n) = \frac{\beta_1 + \beta_2}{2}n + O(1),$$

$$\text{var}(Y_n) = \frac{(\beta_1 - \beta_2)^2}{12}n^2 + O(n)$$

$$\frac{Y_n}{n} \xrightarrow{d} U \quad (\text{convergence in dist. to a uniform r.v.})$$

$$\text{where } f_U(x) = \begin{cases} 0 & \text{if } x \leq b_1 \\ \frac{1}{b_2 - b_1} & \text{if } b_1 < x \leq b_2 \\ 0 & \text{if } b_2 < x \end{cases}$$

$$\text{with } b_1 = \min\{\beta_1, \beta_2\}, b_2 = \max\{\beta_1, \beta_2\}$$

Local limit law of uniform type

[GLV19]

$$f_U(x) = \begin{cases} 0 & \text{if } x \leq b_1 \\ \frac{1}{b_2 - b_1} & \text{if } b_1 < x \leq b_2 \\ 0 & \text{if } b_2 < x \end{cases}, \quad \text{where } \begin{cases} b_1 = \min\{\beta_1, \beta_2\}, \\ b_2 = \max\{\beta_1, \beta_2\} \end{cases}$$

Theorem

Under the same hps $\begin{cases} \lambda_1 = \lambda_2, \beta_1 \neq \beta_2 \\ M_1, M_2 \text{ primitive, } A_j \neq 0 \neq B_j, j = 1, 2 \end{cases}$
 if $(A_1, B_1), (A_2, B_2)$ aperiodic then

$$|n \Pr(Y_n = k) - f_U(x)| = O\left(\frac{(\log n)^{3/2} \tau_n}{\sqrt{n}}\right)$$

for all $k = k(n) \in [n]$ s.t. $k/n \rightarrow x$, $x \in \mathbb{R}$, $\beta_1 \neq x \neq \beta_2$;
 where $\{\tau_n\} \subset \mathbb{R} : \tau_n \rightarrow +\infty, \tau_n = o(\log \log n)$ (arbitrarily slow)

Example

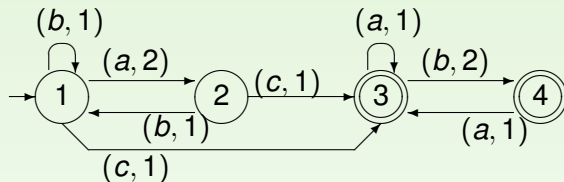


Figure: $\lambda_1 = \lambda_2 = 2$, $1/3 = \beta_1 \neq \beta_2 = 2/3$ (equipotent)

$$L = \{x \in \{a, b\}^* \mid aa \notin x\} \quad c \quad \{y \in \{a, b\}^* \mid bb \notin y\}$$

$$M_1 = M_2, \quad A_1 \neq A_2, \quad (A_1, B_1), (A_2, B_2) \text{ aperiodic}$$

$\Rightarrow Y_n/n$ has local limit of **uniform** type

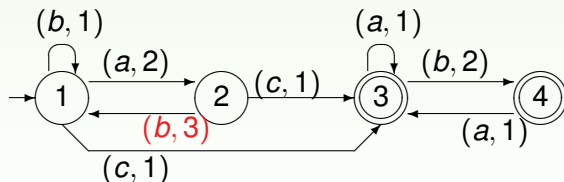


Figure: $\lambda_1 = 3, \lambda_2 = 2 \Rightarrow Y_n/n$ has local limit of **Gaussian** type

Proof (outline)

$$p_n(k) := \Pr(Y_n = k), \quad k = 0, 1, \dots, n$$

Characteristic function

$$\Psi_n(t) = \sum_{j=0}^n p_n(j) e^{itj}$$

Inversion formula

$$p_n(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \Psi_n(t) e^{-itk} dt$$

Ingredients: a) **Saddle Point Method**
 b) analysis of $\Psi_n(t)$ in $[-\pi, \pi]$

a) Split the interval $[-\pi, \pi]$ into 3 sets:

$$\underbrace{\{t \in \mathbb{R} : c < |t| \leq \pi\}}_{(I)}, \quad \underbrace{\{t \in \mathbb{R} : n^{-q} < |t| \leq c\}}_{(II)}, \quad \underbrace{[-n^{-q}, n^{-q}]}_{(III)}$$

for suitable $0 < c < 1, \frac{1}{3} < q < \frac{1}{2}$.

Proof (outline)

(I) From aperiodicity condition

$$\int_{c < |t| \leq \pi} \Psi_n(t) e^{-itk} dt = O(\varepsilon^n) \quad 0 \leq \varepsilon < 1$$

(II)

$$\int_{n^{-q} < |t| \leq c} \Psi_n(t) e^{-itk} dt = O(e^{-an^{1-2q}}) \quad a > 0, \frac{1}{3} < q < \frac{1}{2}$$

(III)

$$\int_{|t| \leq n^{-q}} \Psi_n(t) e^{-itk} dt = \frac{2\pi}{n} f_U(x) + o\left(\frac{1}{n}\right)$$

f_U : PDF of $U \sim \mathcal{U}(\beta_1, \beta_2)$

Of interest here: details of **claim (III)**

$$\Psi_n(t) = \frac{e^{i\beta_1 tn - \frac{\gamma_1}{2} t^2 n} - e^{i\beta_2 tn - \frac{\gamma_2}{2} t^2 n}}{it(\beta_1 - \beta_2)n} + \dots \quad |t| \leq n^{-q}$$

Equipotent bicomponent models with equal β 's and different γ 's

Hp: $\lambda_1 = \lambda_2$, $\beta_1 = \beta_2$, $\gamma_1 \neq \gamma_2$,

M_1, M_2 primitive, $A_j \neq 0 \neq B_j$ ($j = 1, 2$)

Set: $\beta = \beta_1 = \beta_2$, $\gamma = \frac{\gamma_1 + \gamma_2}{2}$

r.v. T mixture of $\mathcal{N}(0, s)$ with

variance s uniformly distributed between $\frac{\gamma_1}{\gamma}$ and $\frac{\gamma_2}{\gamma}$.

$$f_T(x) = \frac{\gamma}{\gamma_2 - \gamma_1} \int_{\frac{\gamma_1}{\gamma}}^{\frac{\gamma_2}{\gamma}} \frac{e^{-\frac{x^2}{2s}}}{\sqrt{2\pi s}} ds \quad \forall x \in \mathbb{R}$$

$f_T(x) \rightsquigarrow$ **heat** equation in dimension 1

It is known that

$$\frac{Y_n - \beta n}{\sqrt{\gamma n}} \xrightarrow{d} T$$

Local limit law of **T** type

[GLV19]

Theorem

Under the same hps $\left\{ \begin{array}{l} \lambda_1 = \lambda_2, \beta_1 = \beta_2, \gamma_1 \neq \gamma_2 \\ M_1, M_2 \text{ primitive, } A_j \neq 0 \neq B_j, j = 1, 2 \end{array} \right.$
 if $(A_1, B_1), (A_2, B_2)$ **aperiodic** and $\gamma = \frac{\gamma_1 + \gamma_2}{2}$ then

$$\left| \sqrt{\gamma n} \Pr(Y_n = k) - f_T \left(\frac{k - \beta n}{\sqrt{\gamma n}} \right) \right| = O \left(\frac{(\log n)^2 \tau_n}{\sqrt{n}} \right)$$

uniformly for $k \in \{0, 1, \dots, n\}$

where $\{\tau_n\} \subset \mathbb{R} : \tau_n \rightarrow +\infty, \tau_n = o(\log \log n)$ (arbitrarily slow)

Equipotent bicomponent models with equal β 's and γ 's

Theorem

Assume the bicomponent model,

let $\lambda_1 = \lambda_2$, $\beta_1 = \beta_2$, $\gamma_1 = \gamma_2$

M_1, M_2 primitive, $A_j \neq 0 \neq B_j$ for $j = 1, 2$.

If $(A_1, B_1), (A_2, B_2)$ are **aperiodic** then

$$\left| \sqrt{n} \Pr(Y_n = k) - \frac{e^{-\frac{(k-\beta n)^2}{2\gamma n}}}{\sqrt{2\pi\gamma}} \right| = O(n^{-1/2})$$

uniformly for every $k \in \{0, 1, \dots, n\}$.

Summary of results in [GLV19]

	Primitive Models	Bicomponent Models			
		$\lambda_1 \neq \lambda_2$	$\lambda_1 = \lambda_2$		
			$\beta_1 \neq \beta_2$	$\beta_1 = \beta_2$ $\gamma_1 \neq \gamma_2$	$\beta_1 = \beta_2$ $\gamma_1 = \gamma_2$
Local limit distribution	$N_{0,1}$ [BCGL03]	$N_{0,1}$	U_{β_1, β_2} [GLV18]	\mathcal{T}	$N_{0,1}$
Convergence rate	$O(n^{-1/2})$	$O(n^{-1/2})$	$O\left(\frac{\tau_n \log^{3/2} n}{\sqrt{n}}\right)$	$O\left(\frac{\tau_n \log^2 n}{\sqrt{n}}\right)$	$O(n^{-1/2})$

Conclusions

- ▶ The results strengthen under aperiodicity conditions the convergence **in distribution** obtained in [dGL04]
- ▶ The convergence rate is $O(n^{-1/2})$ for all Gaussian limits
- ▶ The convergence rate is arbitrarily “**slower**” than $O(n^{-1/2})$ in the other cases
- ▶ For bicomponent models without communication with aperiodicity conditions

$A_0 + B_0 = 0$, $r = s + t$ for $s, t \in \mathbb{R}_+^{\text{rat}} \langle\langle a, b \rangle\rangle$
 we get [GLV19b] (with similar convergence rate):

- ▶ **Gaussian** local limit for Y_n in dominant cases $\lambda_1 \neq \lambda_2$
- ▶ **shifted Bernoullian** local limit if $\lambda_1 = \lambda_2$, $\beta_1 \neq \beta_2$
- ▶ local limit towards a **convex combination of 2 Gaussian laws** if $\lambda_1 = \lambda_2$, $\beta_1 = \beta_2$ and $\gamma_1 \neq \gamma_2$
- ▶ **Gaussian** local limit for Y_n if $\lambda_1 = \lambda_2$, $\beta_1 = \beta_2$, $\gamma_1 = \gamma_2$

Thank you !

Theorem [Nicodeme-Salvy-Flajolet 02, BCGL03]

For every regular $R \subseteq \Sigma^*$ (pattern) and every Markovian source (π, P) over Σ there exists $r \in \mathbb{R}_+^{\text{rat}} \langle\langle a, b \rangle\rangle$ s.t.

$$\Pr(O_n(\pi, P, R) = k) = \Pr(Y_n(r) = k) \quad \forall k = 0, \dots, n$$

[Proof sketch]

Consider a det. f.s. automaton (Q, p, δ, F) recognizing $\Sigma^* R$.

Define a new set of states $Q' = \{p\} \cup \{(q, \sigma) \mid q \in Q, \sigma \in \Sigma\}$.

Define a linear representation (ξ, A, B, η) over Q' such that:

► $\xi = \chi_p, \eta = \chi_{Q'}$

► $A_{p,(q,\sigma)} = \begin{cases} \pi_\sigma & \text{if } q = \delta(p, \sigma) \wedge q \in F \\ 0 & \text{otherwise} \end{cases},$

$$A_{(q,\sigma),(q',\tau)} = \begin{cases} P(\sigma, \tau) & \text{if } q' = \delta(q, \sigma) \wedge q' \in F \\ 0 & \text{otherwise} \end{cases}, \quad \forall \sigma, \tau \in \Sigma, q, q' \in Q$$

► $B_{p,(q,\sigma)} = \begin{cases} \pi_\sigma & \text{if } q = \delta(p, \sigma) \wedge q \notin F \\ 0 & \text{otherwise} \end{cases},$

$$B_{(q,\sigma),(q',\tau)} = \begin{cases} P(\sigma, \tau) & \text{if } q' = \delta(q, \sigma) \wedge q' \notin F \\ 0 & \text{otherwise} \end{cases}, \quad \forall \sigma, \tau \in \Sigma, q, q' \in Q$$

(Continue)

Let $\mu : \{a, b\}^* \rightarrow \mathbb{R}_+^{Q' \times Q'}$ be the morphism generated by $A = \mu(a)$, $B = \mu(b)$. Then, for every $n \in \mathbb{N}$ and every $k \in \{0, 1, \dots, n\}$, we have

$$\Pr(Y_n = k) = \frac{\sum_{x \in \{a, b\}^n, |x|_a = k} \xi' \mu(x) \eta}{\xi'(A + B)^n \eta} = \frac{\sum_{w \in \Sigma^n, |w|_R = k} \Pr(w)}{1} = \Pr(O_n = k)$$

Note that $A + B$ is a stochastic matrix and the bivariate generating function

$$\sum_{n=0}^{+\infty} \sum_{k=0}^n \Pr(O_n = k) x^k y^n = \xi' (I - (Ax + By))^{-1} \eta$$

is rational.

Theorem [BCGL03]

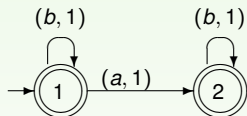
For some $r \in \mathbb{R}_+^{\text{rat}} \langle\langle a, b \rangle\rangle$ the generating function

$$f_r(x, y) = \sum_{n=0}^{+\infty} \sum_{k=0}^n \Pr(Y_n = k) x^k y^n$$

is not rational.

[Proof sketch]

Consider the series r defined by the weighted automaton



$$r(x) = \begin{cases} 1 & \text{if } |x|_a \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Then $\Pr(Y_n = k) = \begin{cases} \frac{1}{n+1} & \text{if } k = 0 \\ \frac{n}{n+1} & \text{if } k = 1 \\ 0 & \text{otherwise} \end{cases}$, and

$$\sum_{n=0}^{+\infty} \sum_{k=0}^n \Pr(Y_n = k) x^k y^n = (x-1) \frac{\log(1-z)}{y} + \frac{x}{1-y}$$