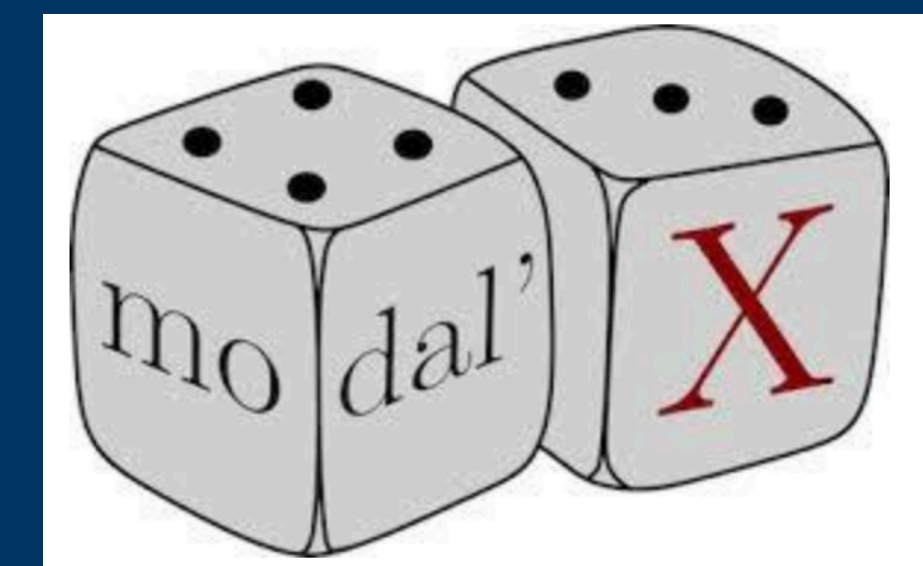


Malliavin calculus for marked binomial processes Chen-Stein method

Hélène Halconruy



Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$.

000110010101011111100001010011011111111111000001101111111110010

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

00011001010101011111010000101001101111111111100000110111111010010

$m_n=5$ $m_n=5$ $m_n=5$ $m_n=5$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\boxed{\{R_n < m_n\} = \{U_n = 0\}}$$

U_n : nb of starting positions
in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 38, 52$.

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

00011001010101011111010000101001101111111111100000110111111010010

$m_n=5$ $m_n=5$ $m_n=5$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\{R_n < m_n\} = \{U_n = 0\}$$

U_n : nb of starting positions in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 34, 38, 52, 53$.

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

0001100101010101111101000010100110111111111100000110111111010010

$\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\boxed{\{R_n < m_n\} = \{U_n = 0\}}$$

U_n : nb of starting positions
in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 34, 35, 36, 37, 38, 39, 52, 53$

Clumps

$$\mathbf{P}(\text{run at } j + 1 \mid \text{run at } j) = p$$

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

0001100101010101111101000010100110111111111100000110111111010010

$\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\boxed{\{R_n < m_n\} = \{U_n = 0\}}$$

U_n : nb of starting positions
in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 34, 35, 36, 37, 38, 39, 52, 53$

Clumps $\mathbf{P}(\text{run at } j+2 \mid \text{run at } j) = p^2$

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

000110010101010111110100001010011011111111110000011011111010010

$\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\boxed{\{R_n < m_n\} = \{U_n = 0\}}$$

U_n : nb of starting positions in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 34, 35, 36, 37, 38, 39, 52, 53$

Clumps $\mathbf{P}(\text{run at } j+k \mid \text{run at } j) = p^k$

Introduction

Application 1 : The longest head run

Large number of independent tosses of a coin, $C_i \sim \mathcal{B}(p)$. m_n : test length

0001100101010101111101000010100110111111111100000110111111010010

$\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$ $\underbrace{\hspace{10em}}_{m_n=5}$

Question : Law of R_n longest series of 1's starting in the first n tosses = *longest head run* ?

$$\boxed{\{R_n < m_n\} = \{U_n = 0\}}$$

U_n : nb of starting positions
in $\{1, \dots, n\}$ of a head run

Runs of length m_n at positions $j = 14, 33, 34, 35, 36, 37, 38, 39, 52, 53$

Clumps Average size of clump : $1 + p + p^2 + \dots$

Application 2 : a rare word in a DNA sequence

Series $X_1X_2\dots X_n$ from $\mathcal{A} = \{A,C,G,T\}$

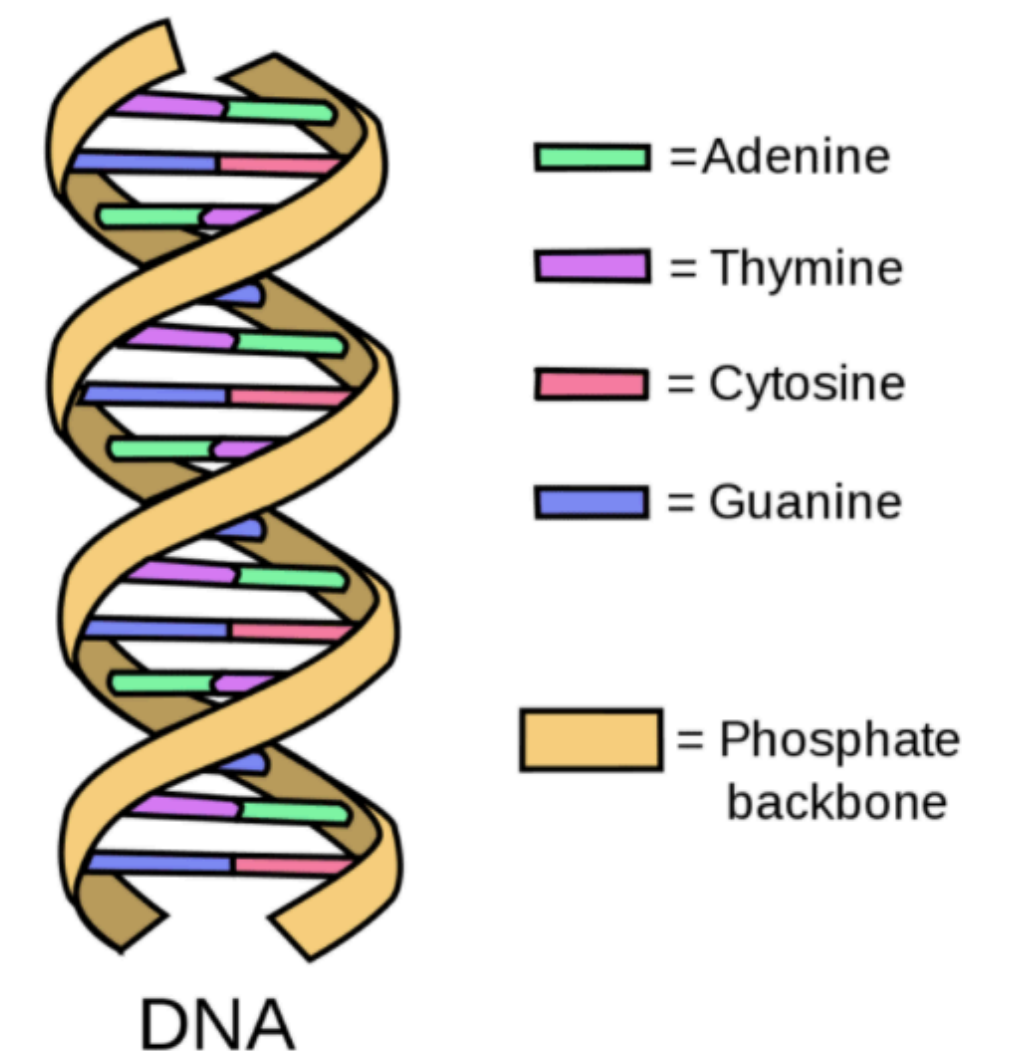
Rare word : $W_n = w_1w_2\dots w_{h_n}$ of length $h_n \sim \log(n)$

Question : number of occurrences of W_n ?

$$\mathfrak{Z}(W_n) = \sum_{j \in I} Z_j = \sum_{j \in I} \mathbf{1}_{\{X_j=w_1, \dots, X_{j+h_n-1}=w_{h_n}\}},$$

Example : $W_n = \text{ACTAA}$.

GACTAACTAAACTAATGAAACTAACG



$$\mathbf{E}[Z_j] = \mu(W_n)$$

$$I = \{1, \dots, n - h_n + 1\}$$

Application 2 : a rare word in a DNA sequence

Series $X_1X_2\dots X_n$ from $\mathcal{A} = \{A,C,G,T\}$

Rare word : $W_n = w_1w_2\dots w_{h_n}$ of length $h_n \sim \log(n)$

Question : number of occurrences of W_n ?

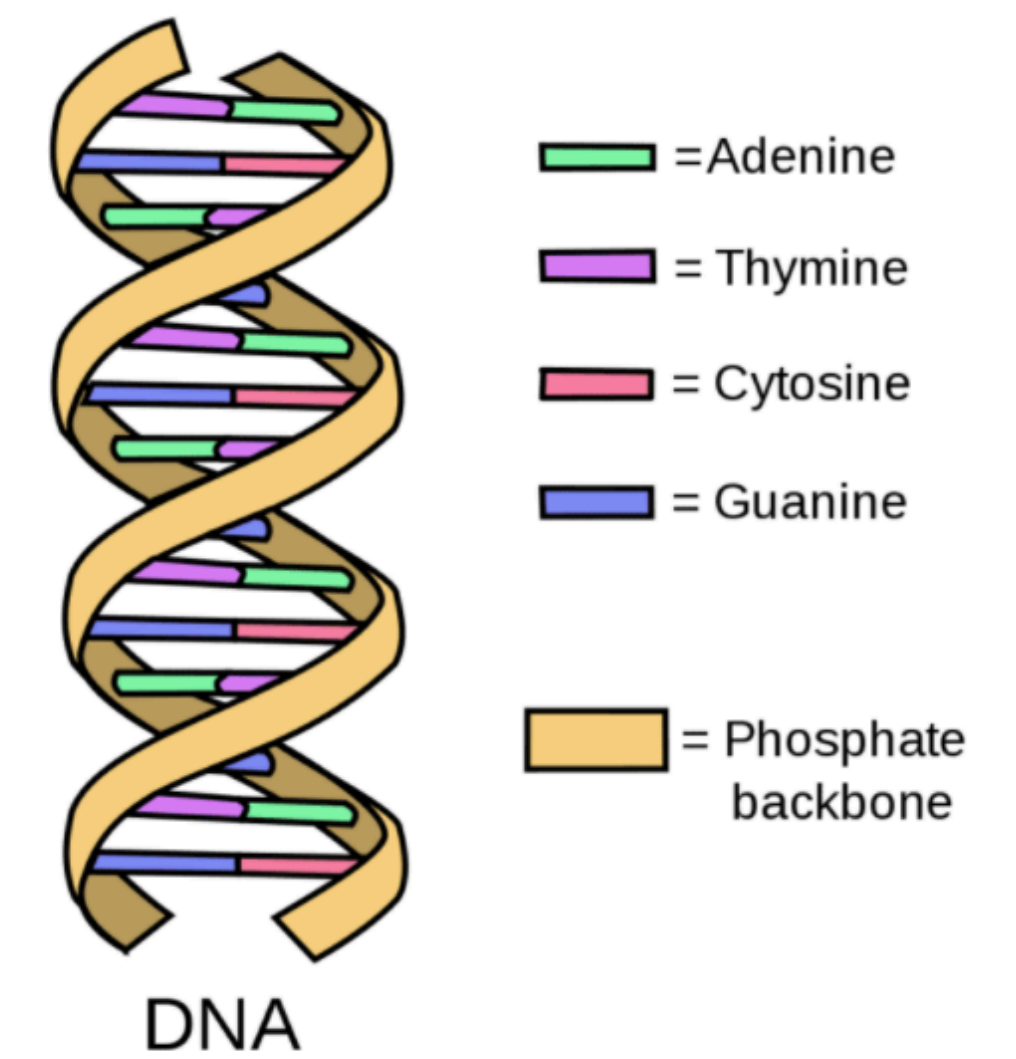
$$\mathfrak{Z}(W_n) = \sum_{j \in I} Z_j = \sum_{j \in I} \mathbf{1}_{\{X_j=w_1, \dots, X_{j+h_n-1}=w_{h_n}\}},$$

Example : $W_n = \text{ACTAA}$.

GACTAACTAAACTAATGAAACTAACG
 ACTAA
 GACTAACTAAACTAATGAAACTAACG

Clumps

3-clump at $j = 2$, 1-clump at $j = 20$



$$\mathbf{E}[Z_j] = \mu(W_n)$$

$$I = \{1, \dots, n - h_n + 1\}$$

Motivation Poisson approximation ?

Average size of clump : $1 + p + p^2 + \dots$

Head runs

Declump

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, C_i \sim \text{Ber}(p),$$

$$U_n = \sum_{k \in I} X_k^n : \text{nb of runs of length } m_n$$

DNA sequence

Motivation Poisson approximation ?

Average size of clump : $1 + p + p^2 + \dots$

Head runs

$$U_n \sim \mathcal{P}(\lambda_n) ?$$

Declump

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, C_i \sim \text{Ber}(p),$$

$$U_n = \sum_{k \in I} X_k^n : \text{nb of runs of length } m_n$$

DNA sequence

Motivation Poisson approximation ?

Average size of clump : $1 + p + p^2 + \dots$

Head runs

$$U_n \sim \mathcal{P}(\lambda_n) ?$$

Declump

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, C_i \sim \text{Ber}(p),$$

$$U_n = \sum_{k \in I} X_k^n : \text{nb of runs of length } m_n$$

DNA sequence

$$\mathbf{P}(\text{occ . word once at } j) = \lambda_n$$

Motivation Poisson approximation ?

Average size of clump : $1 + p + p^2 + \dots$

Head runs

$$U_n \sim \mathcal{P}(\lambda_n) ?$$

Declump

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, C_i \sim \text{Ber}(p),$$

$$U_n = \sum_{k \in I} X_k^n : \text{nb of runs of length } m_n$$

DNA sequence

$$\mathfrak{L}(W_n) \sim \mathcal{PC}(\lambda_n, \mathbf{V}) ?$$

$$\mathbf{P}(\text{occ. word } k \text{ times at } j) = \lambda_n \mathbf{V}(\{k\})$$

Motivation Poisson approximation ?

Average size of clump : $1 + p + p^2 + \dots$

Head runs

$$U_n \sim \mathcal{P}(\lambda_n) ?$$

Declump

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, C_i \sim \text{Ber}(p),$$

$$U_n = \sum_{k \in I} X_k^n : \text{nb of runs of length } m_n$$

DNA sequence

$$\mathfrak{L}(W_n) \sim \mathcal{PC}(\lambda_n, \mathbf{V}) ?$$

$$\mathbf{P}(\text{occ. word } k \text{ times at } j) = \lambda_n \mathbf{V}(\{k\})$$

Question : How to quantify the convergence rate for TV- distance ?

For \mathbf{P}, \mathbf{Q} probability measures on $\mathbb{N} = \{0, 1, \dots\}$ $d_{\text{TV}}(\mathbf{P}, \mathbf{Q}) = \sup_{A \subseteq \mathbb{N}} |\mathbf{P}(A) - \mathbf{Q}(A)|$

Chen-Stein method

Equivalent problem $\mathbf{P}^* = \mathcal{PC}(\lambda_0, \mathbf{V}), F : \mathbb{N}\text{-valued}$

L : Stein operator $L\varphi(F) = \lambda_0 \sum_{k \geq 1} k\varphi(F+k)\mathbf{V}(\{k\}) - F\varphi(F)$

$(\psi_A, A \subset \mathbb{N})$ Stein's class

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) = \sup_{A \subset \mathbb{N}} \left| \mathbf{E} \left[\lambda_0 \sum_{k \geq 1} k\psi_A(F+k)\mathbf{V}(\{k\}) - F\psi_A(F) \right] \right|$$

Original approach

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) = \sup_{A \subset \mathbb{N}} \left| \mathbf{E} \left[\lambda_0 \sum_{k \geq 1} k \psi_A(F + k) \mathbf{V}(\{k\}) - F \psi_A(F) \right] \right|$$

Frame $F = \sum_{i \in I} \sum_{k \geq 1} X_i \mathbf{1}_{\{\text{occ. at } i \text{ of size } k\}}, \quad X_i = \mathbf{1}_{\{\text{occ. at } i\}} \sim \text{Ber}(p), \quad I \subset \mathbb{N}.$

Neighborhoods of dependence $(\mathcal{N}_i, i \in I)$ X_i is independent of $\{X_j, j \notin \mathcal{N}_i \cup \{i\}\}$

Approximation

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) \leq b_1 + b_2 + b_3$$

$$b_1 = \sum_{i \in I} (\mathbf{E}[X_i])^2 + \sum_{i \in I} \sum_{j \in \mathcal{N}_i, j \neq i} (\mathbf{E}[X_i])(\mathbf{E}[X_j])$$

$$b_2 = \sum_{i \in I} \sum_{j \in \mathcal{N}_i} \mathbf{E}[X_i X_j]$$

$$b_3 = \sum_{i \in I} \mathbf{E} \left[\left| \mathbf{E}[X_i - p_i | \sigma(X_j, j \notin \mathcal{N}_i)] \right| \right]$$

New approach

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) = \sup_{A \subseteq \mathbb{N}} \left| \mathbf{E} \left[\lambda_0 \sum_{k \geq 1} k \psi_A(F + k) \mathbf{V}(\{k\}) - F \psi_A(F) \right] \right|$$

Idea 1.

$$\begin{aligned} F &= \sum_{i \in I} \sum_{k \geq 1} X_i \mathbf{1}_{\{\text{occ. at } i \text{ of size } k\}} \\ &= f((\Delta N_1, V_1), \dots, (\Delta N_n, V_n), \dots) \end{aligned}$$

$$(\Delta N_j)_j \perp\!\!\!\perp (V_j)_j.$$

$$\Delta N_j \sim \mathcal{B}(\lambda), V_j \sim \mathbf{V}$$

New approach

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) = \sup_{A \subseteq \mathbb{N}} \left| \mathbf{E} \left[\lambda_0 \sum_{k \geq 1} k \psi_A(F + k) \mathbf{V}(\{k\}) - F \psi_A(F) \right] \right|$$

Idea 1.

$$\begin{aligned} F &= \sum_{i \in I} \sum_{k \geq 1} X_i \mathbf{1}_{\{\text{occ. at } i \text{ of size } k\}} \\ &= f((\Delta N_1, V_1), \dots, (\Delta N_n, V_n), \dots) \end{aligned}$$

$$(\Delta N_j)_j \perp\!\!\!\perp (V_j)_j.$$

$$\Delta N_j \sim \mathcal{B}(\lambda), V_j \sim \mathbf{V}$$

Idea 2.

Nourdin Peccati 2009

New approach

$$d_{\text{TV}}(\mathbf{P}_F, \mathbf{P}^*) = \sup_{A \subseteq \mathbb{N}} \left| \mathbf{E} \left[\lambda_0 \sum_{k \geq 1} k \psi_A(F + k) \mathbf{V}(\{k\}) - F \psi_A(F) \right] \right|$$

Idea 1.

$$F = \sum_{i \in I} \sum_{k \geq 1} X_i \mathbf{1}_{\{\text{occ. at } i \text{ of size } k\}}$$
$$= f((\Delta N_1, V_1), \dots, (\Delta N_n, V_n), \dots)$$

$$(\Delta N_j)_j \perp\!\!\!\perp (V_j)_j.$$

$$\Delta N_j \sim \mathcal{B}(\lambda), V_j \sim \mathbf{V}$$

Idea 2.

$\mathbf{E}[F \psi_A(F)]$

Nourdin Peccati 2009

Integration by parts ?

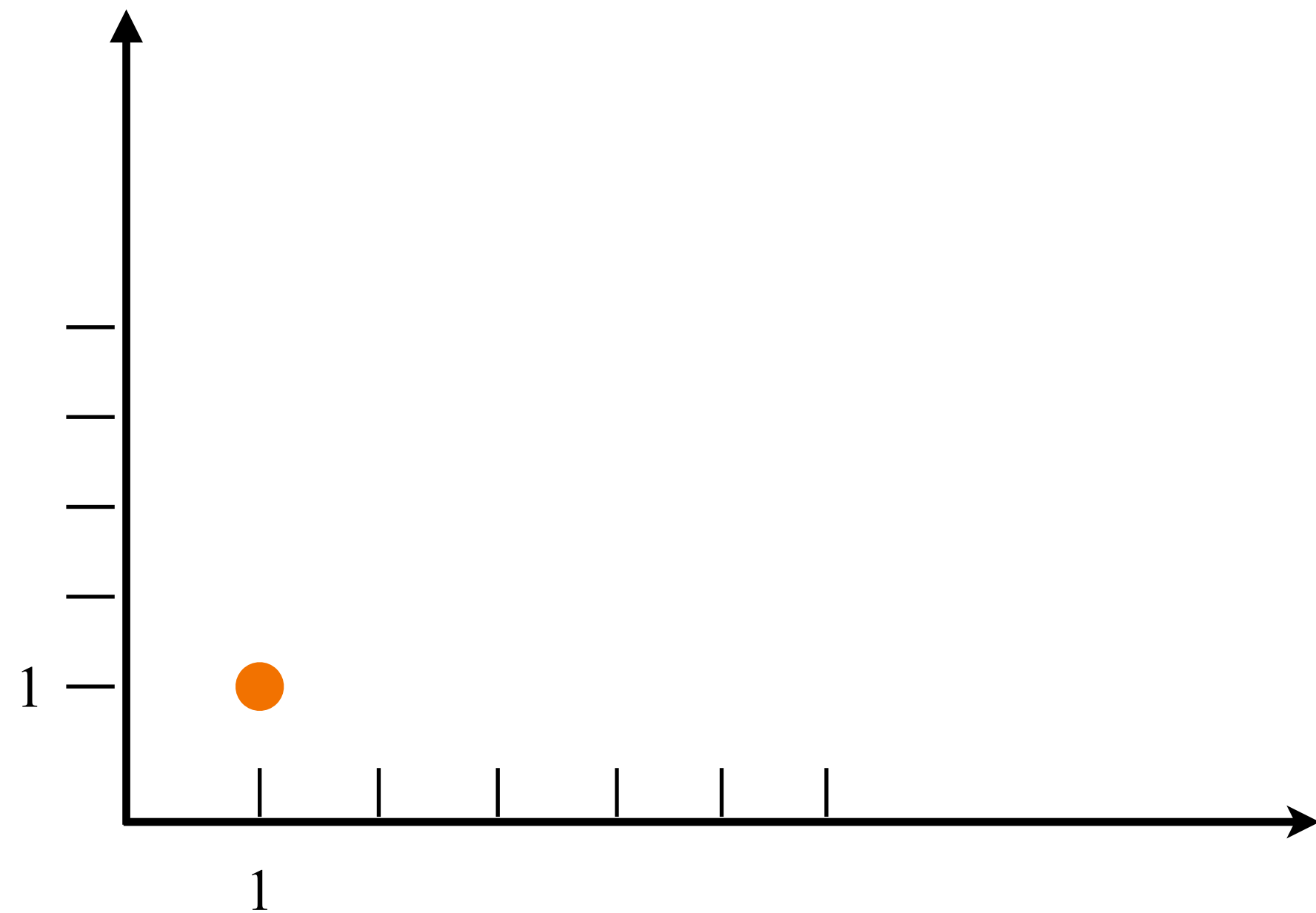
Malliavin calculus for marked binomial processes (MBP)

Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

E countable = *mark space*.

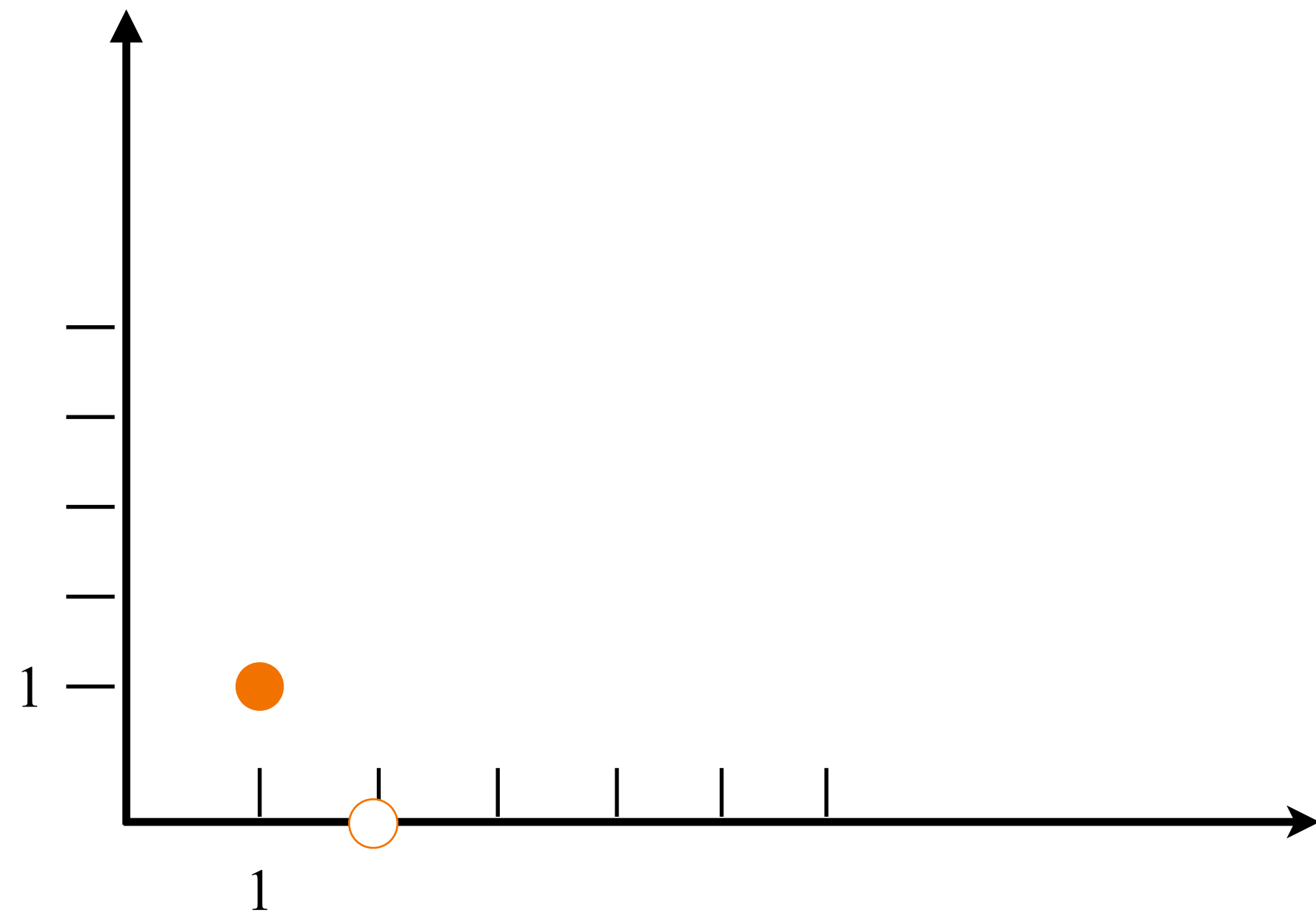
$$\underline{n = 1} \quad \Delta N_1 = 1, V_1 = 1$$
$$\eta(1,1) = 1$$



Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

E countable = *mark space*.



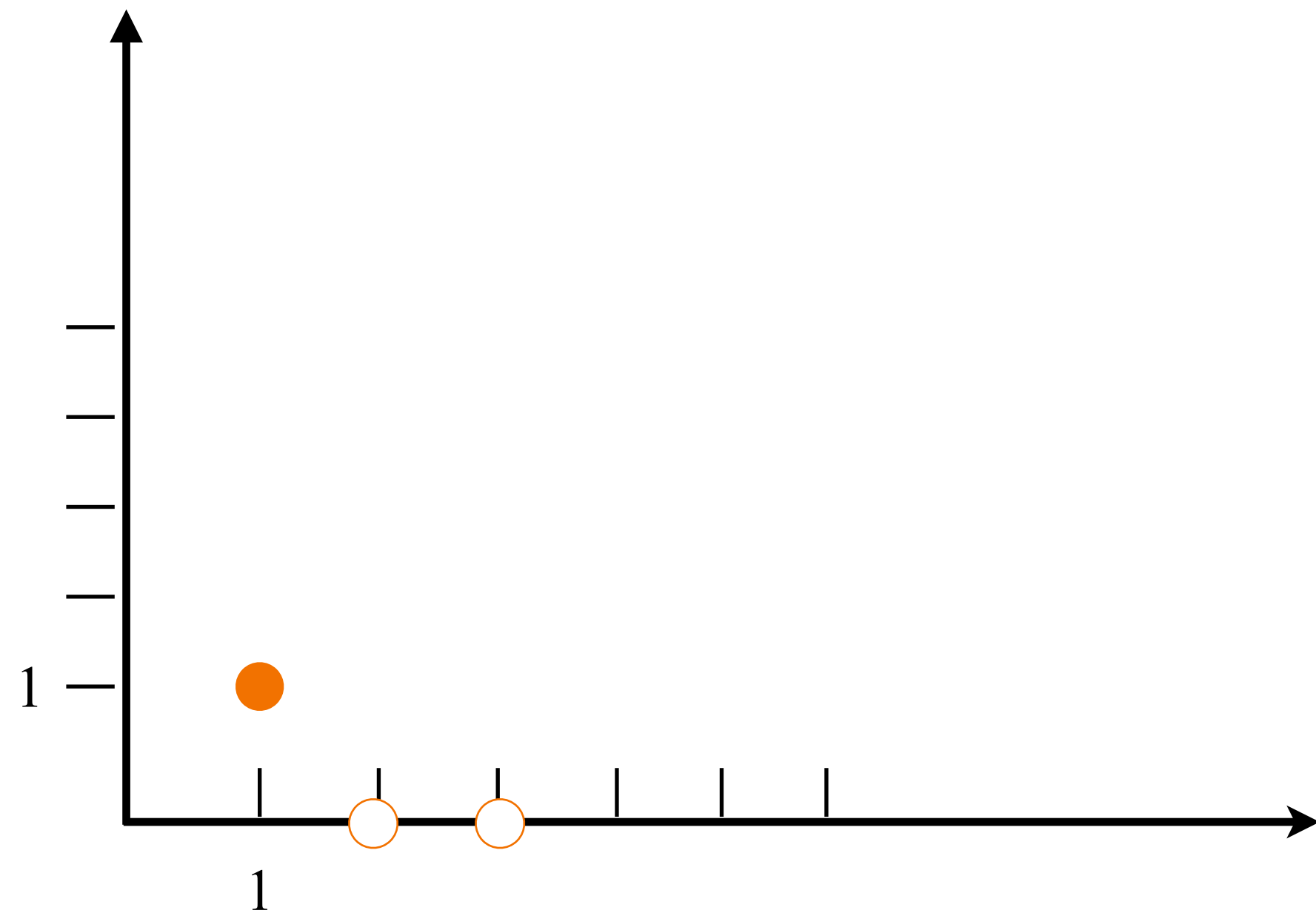
$$\underline{n = 1} \quad \Delta N_1 = 1, V_1 = 1$$
$$\eta(1,1) = 1$$

$$\underline{n = 2} \quad \Delta N_2 = 0$$
$$\eta(2,\cdot) = \sum_{k \in E} \eta(2,k) = 0$$

Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

E countable = *mark space*.



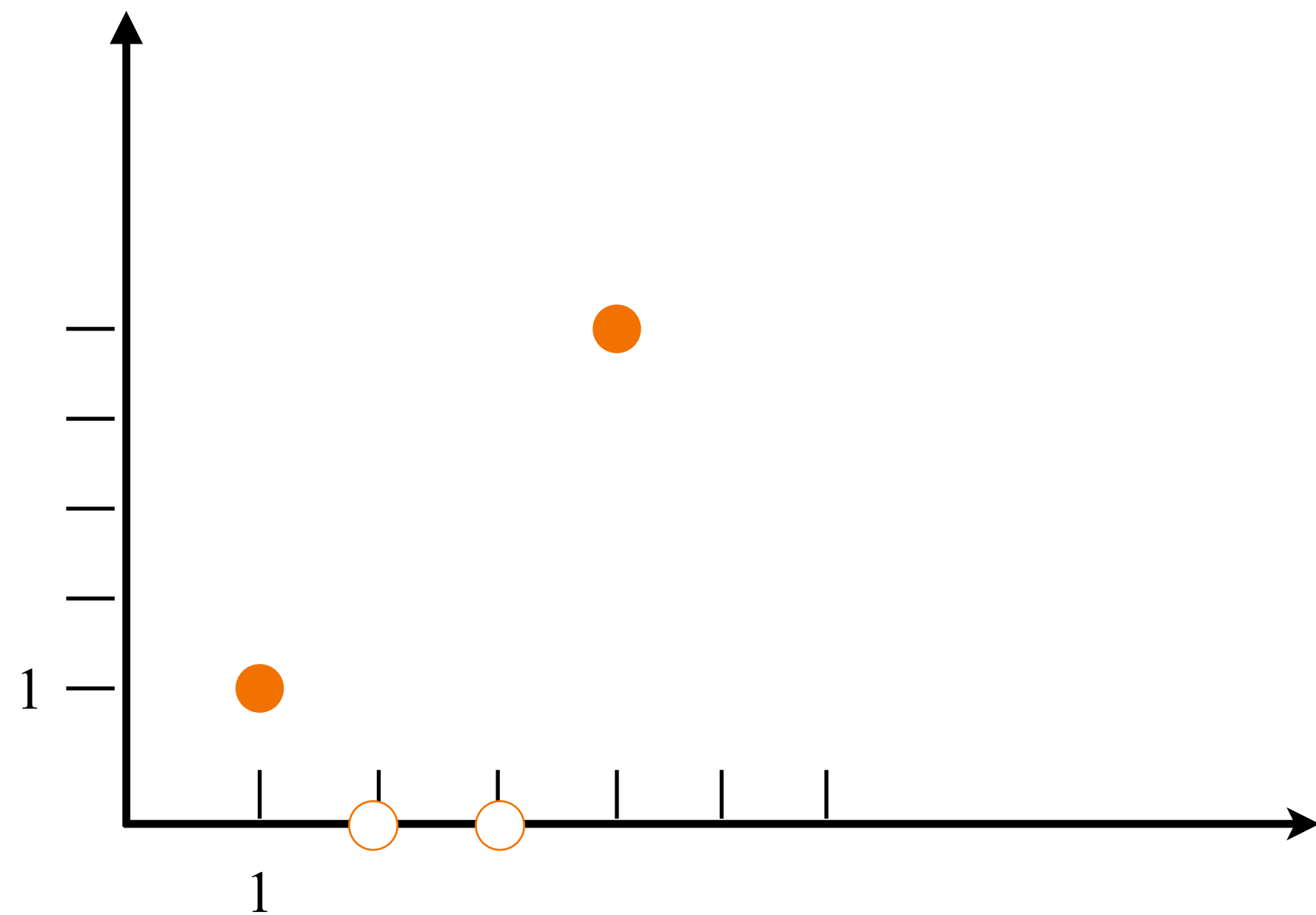
$$\underline{n = 1} \quad \Delta N_1 = 1, V_1 = 1$$
$$\eta(1,1) = 1$$

$$\underline{n = 2,3} \quad \Delta N_2, \Delta N_3 = 0$$
$$\eta(2,\cdot) = \sum_{k \in E} \eta(2,k) = 0$$

Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

E countable = *mark space*.



$$\underline{n = 1} \quad \Delta N_1 = 1, V_1 = 1$$
$$\eta(1,1) = 1$$

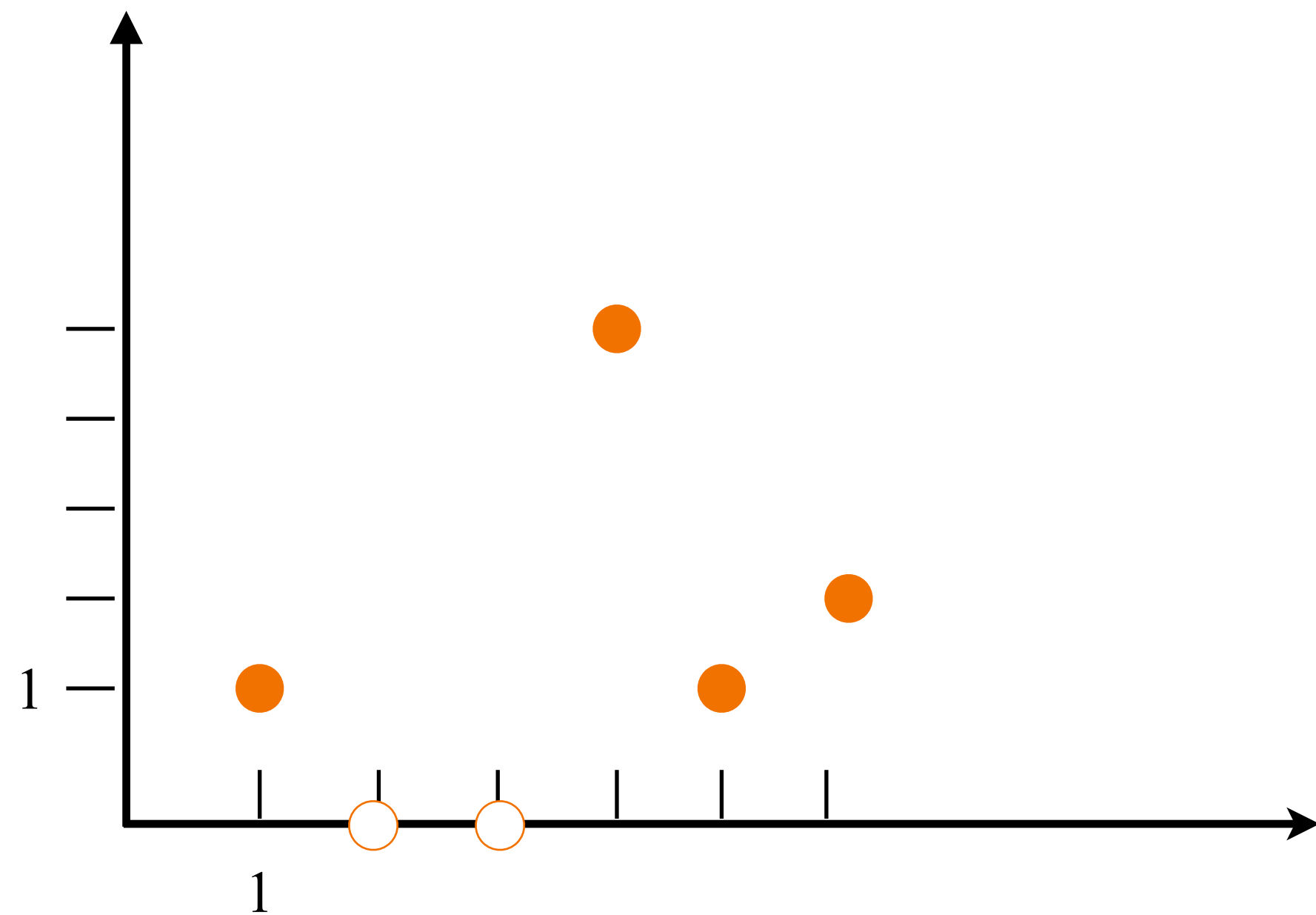
$$\underline{n = 2,3} \quad \Delta N_2, \Delta N_3 = 0$$
$$\eta(2,\cdot) = \sum_{k \in E} \eta(2,k) = 0$$

$$\underline{n = 4} \quad \Delta N_4 = 1, V_2 = 5$$
$$\eta(4,5) = 1$$

Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

E countable = *mark space*.



$$\underline{n = 1} \quad \Delta N_1 = 1, V_1 = 1$$

$$\eta(1,1) = 1$$

$$\underline{n = 2,3} \quad \Delta N_2, \Delta N_3 = 0$$

$$\eta(2,\cdot) = \sum_{k \in E} \eta(2,k) = 0$$

$$\underline{n = 4} \quad \Delta N_4 = 1, V_2 = 5$$

$$\eta(4,5) = 1$$

$$\underline{n = 5,6} \quad \text{Etc.}$$

Framework

Marked binomial processes $(\Omega, \mathcal{F}, \mathbf{P})$

MBP (λ, \mathbf{V}) $\mathbb{X} = \mathbb{N}^* \times \mathbf{E}$

$$\eta = \sum_{t=1}^{\infty} \mathbf{1}_{\{T_t < \infty\}} \delta_{(T_t, V_t)}$$

T_t t -th jump, $T_0 = 0$ and $T_t - T_{t-1} \sim \mathcal{G}(\lambda)$, i.i.d.

$$N_0 = 0, \Delta N_t = N_t - N_{t-1} = \mathbf{1}_{\{\text{jump at } t\}}$$

V_t mark of the t -th jump i.i.d. $V_t \sim \mathbf{V}$, $(V_t)_t \perp\!\!\!\perp (T_t)_t$

Intensity of MBP (λ, \mathbf{V})

$$d\nu(t, k) = \lambda \mathbf{V}(\{k\}) \delta_{(t,k)}$$

$$\Delta N_t \sim \text{Ber}(\lambda)$$

Functionals and processes $(\Omega, \mathcal{A}, \mathbf{P})$

Marked binomial functionals

$$F = \mathfrak{f}(\eta) \quad \mathfrak{f} \text{ representative of } F$$

$$L^p(\mathbf{P}) = \{F : \mathbf{E}[|F|^p] < \infty\} \quad p \geq 1$$

Marked binomial processes

$$u = \sum_{(t,k) \in \mathbb{X}} \mathbf{u}(\eta, (t, k)) \mathbf{1}_{(t,k)} \quad \mathbf{u} \text{ representative of } u$$

$$L^p(\mathbf{P} \otimes \nu) = \left\{ u : \mathbf{E} \left[\sum_{(t,k)} |u(\cdot, (t, k))|^p \right] < \infty \right\}$$

Filtration $(\mathcal{F}_t)_{t \geq 1}$ $\mathcal{F}_t = \sigma\{\eta(s, k), s \leq t, k \in E\}.$

L^1 -theory

Mecke formula

\mathbf{u} representative
of u

Lemma (H)

For \mathbf{u} non-negative,

$$\mathbf{E} \left[\sum_{(t,k) \in \eta} \mathbf{u}(\eta, (t, k)) \right] = \mathbf{E} \left[\int_{\mathbb{X}} \mathbf{u}(\pi_t(\eta) + \delta_{(t,k)}, (t, k)) d\nu(t, k) \right]$$

L^1 -theory

Mecke formula

Lemma (H)

For u non-negative,

$$\mathbf{E} \left[\sum_{(t,k) \in \eta} u(\eta, (t, k)) \right] = \mathbf{E} \left[\int_{\mathbb{X}} u(\pi_t(\eta) + \delta_{(t,k)}, (t, k)) d\nu(t, k) \right]$$

where

$$\pi_t(\eta) = \sum_{s \neq t} \sum_{k \in E} \eta(s, k)$$

L^1 -Operators and integration by parts

Add-one cost operator

$$D_{(t,k)}^+ F = \mathfrak{f}(\pi_t(\eta) + \delta_{(t,k)}) - \mathfrak{f}(\pi_t(\eta))$$

Proposition. L^1 integration by parts (H)

For u predictable, F s.t. $\widetilde{D}Fu \in L^1(\mathbf{P} \otimes \nu)$,

$$\mathbf{E} \left[\int_{\mathbb{X}} \widetilde{D}_{(t,k)} F u_{(t,k)} d\nu(t, k) \right] = \mathbf{E}[F \widetilde{\delta}(u)]$$

where $\widetilde{D}_{(t,k)} F = \mathfrak{f}(\pi_t(\eta) + \delta_{(t,k)}) - \mathfrak{f}(\eta)$.

« Skorokhod integral »

$$\widetilde{\delta}(u) = \sum_{(t,k) \in \mathbb{X}} u_{(t,k)} \Delta Z_{(t,k)}$$

$$\Delta Z_{(t,k)} := \mathbf{1}_{\{\eta(s,k)=1\}} - \lambda \mathbf{V}(\{k\})$$

- $\mathbf{E}[\Delta Z_{(t,k)}] = 0$
- $\Delta Z_{(t,\cdot)} \perp\!\!\!\perp \Delta Z_{(s,\cdot)}$ for $t \neq s$
- $\mathbf{E}[\Delta Z_{(t,k)} \Delta Z_{(t,\ell)}] \neq 0$.

L^1 -Operators and integration by parts

Add-one cost operator

$$D_{(t,k)}^+ F = \mathfrak{f}(\pi_t(\eta) + \delta_{(t,k)}) - \mathfrak{f}(\pi_t(\eta))$$

Proposition. L^1 integration by parts (H)

For u predictable, F s.t. $\widetilde{D}Fu \in L^1(\mathbf{P} \otimes \nu)$,

$$\mathbf{E} \left[\int_{\mathbb{X}} \widetilde{D}_{(t,k)} F u_{(t,k)} d\nu(t, k) \right] = \mathbf{E}[F \widetilde{\delta}(u)]$$

where $\widetilde{D}_{(t,k)} F = \mathfrak{f}(\pi_t(\eta) + \delta_{(t,k)}) - \mathfrak{f}(\eta)$.

« Skorokhod integral »

$$\widetilde{\delta}(u) = \sum_{(t,k) \in \mathbb{X}} u_{(t,k)} \Delta Z_{(t,k)}$$

$$\Delta Z_{(t,k)} := \mathbf{1}_{\{\eta(s,k)=1\}} - \lambda \mathbf{V}(\{k\})$$

- $\mathbf{E}[\Delta Z_{(t,k)}] = 0$
- $\Delta Z_{(t,\cdot)} \perp\!\!\!\perp \Delta Z_{(s,\cdot)}$ for $t \neq s$
- $\mathbf{E}[\Delta Z_{(t,k)} \Delta Z_{(t,\ell)}] \neq 0$.

L^2 -theory

The family \mathcal{R}

$$\Delta Z_{(t,k)} = \mathbf{1}_{\{(t,k) \in \eta\}} - \lambda \mathbf{V}(\{k\})$$

$$\Delta R_0 = 1, \quad \Delta R_{(t,1)} = \Delta Z_{(t,1)} \quad \text{and} \quad \Delta R_{(t,k)} = \Delta Z_{(t,k)} - \sum_{j=1}^{k-1} \frac{\mathbf{E}[\Delta Z_{(t,k)} \Delta R_{(t,j)}]}{\mathbf{E}[(\Delta R_{(t,j)})^2]} \Delta R_{(t,j)}$$

The $\Delta R_{(t,k)}$ satisfy

- $\mathbf{E}[\Delta R_{(t,k)}] = 0$
- $\Delta R_{(t,\cdot)} \perp\!\!\!\perp \Delta R_{(s,\cdot)}$ for $t \neq s$
- $\mathbf{E}[\Delta R_{(t,k)} \Delta R_{(t,\ell)}] = 0$ for $k \neq \ell$ and $\mathbf{E}[(\Delta R_{(t,k)})^2] =: \kappa_k$
- $\Delta R_{(t,k)} = \Delta Z_{(t,k)} - \sum_{j=1}^{k-1} \rho_j \Delta Z_{(t,j)}$

Chaotic decomposition $(\mathbf{t}_n, \mathbf{k}_n) = ((t_1, k_1), \dots, (t_n, k_n))$

$$\begin{aligned}\Delta Z_{(t,k)} &= \mathbf{1}_{(t,k)} - \lambda \mathbf{V}(\{k\}) \\ \Delta R_{(t,k)} &= \Delta Z_{(t,k)} - \sum_{j=1}^{k-1} \rho_j \Delta Z_{(t,j)} \\ \mathbf{E}[(\Delta R_{(t,k)})^2] &=: \kappa_k\end{aligned}$$

Hilbert space of symmetric functions

$$\langle f_n, g_n \rangle_{L^2(\tilde{\nu})^{\circ n}} = n! \int_{\mathbb{X}^{n,<}} f_n(\mathbf{t}_n, \mathbf{k}_n) g_n(\mathbf{t}_n, \mathbf{k}_n) d\tilde{\nu}^{\otimes n}(\mathbf{t}_n, \mathbf{k}_n) \text{ with } \tilde{\nu}(t, k) = \kappa_k \nu(t, k)$$

Proposition. \mathcal{R} -integration (H)

The \mathcal{R} -integral of $f_n \in \mathcal{C}^n(\mathbb{X}^n, \mathbb{R})$ is

$$J_0(f_0) = f_0 \quad \text{and} \quad J_n(f_n) = n! \sum_{(\mathbf{t}_n, \mathbf{k}_n)} f_n(\mathbf{t}_n, \mathbf{k}_n) \prod_{i=1}^n \Delta R_{(t_i, k_i)} \quad (n \geq 1)$$

Isometry property : for $f_n \in L^2(\tilde{\nu})^{\circ n}$, $g_m \in L^2(\tilde{\nu})^{\circ m}$,

$$\mathbf{E}[J_n(f_n) J_m(g_m)] = \mathbf{1}_n(m) n! \langle f_n, g_n \rangle_{L^2(\tilde{\nu})^{\circ n}}$$

Chaotic decomposition

$$J_n(f_n) = n! \sum_{(\mathbf{t}_n, \mathbf{k}_n)} f_n(\mathbf{t}_n, \mathbf{k}_n) \prod_{i=1}^n \Delta R_{(t_i, k_i)}$$

Marked binomial functionals

$$F = f_0 \mathbf{1}_{\{\eta(\mathbb{X})=0\}} + \sum_{n \geq 1} \sum_{(\mathbf{t}_n, \mathbf{k}_n) \in \mathbb{X}^n} f_n(\mathbf{t}_n, \mathbf{k}_n) \mathbf{1}_{\{\eta(\mathbb{X})=n\}} \prod_{i=1}^n \mathbf{1}_{\{\eta(t_i, k_i)=1\}} ; f_n \in L^2(\tilde{\mathcal{V}})^{\circ n}$$

Marked binomial chaoses

$$\mathcal{H}_0 = \mathbf{R} \quad \text{and} \quad \mathcal{H}_n = \{J_n(f_n) ; f_n \in L^2(\tilde{\mathcal{V}})^{\circ n}\}$$

Theorem. Chaotic decomposition (H)

$$L^2(\mathbf{P}) = \bigoplus_{n \in \mathbf{N}} \mathcal{H}_n$$

In other words, $F \in L^2(\mathbf{P})$ can be uniquely expanded

$$F = \mathbf{E}[F] + \sum_{n \in \mathbf{N}^*} J_n(f_n)$$

L^2 -Integration by parts - Malliavin calculus

Annihilation operator

$$D_{(t,k)}J_n(f_n) = nJ_{n-1}(f_n(\star, (t, k))\mathbf{1}_{\{1, \dots, t-1\}^n})$$

Itô-Wiener integral

$$\delta(u) = J_1(u) = \sum_{(t,k)} u_{(t,k)} \Delta R_{(t,k)}$$

Proposition. (Extended) L^2 -Integration by parts (H)

For any $(F, u) \in \text{dom}D \times \text{dom}\delta$,

$$\mathbf{E}[F\delta u] = \mathbf{E}[\langle DF, u \rangle_{L^2(\mathbb{X}, \tilde{\nu})}]$$

$L^1 - L^2$ correspondence?

$$\Delta Z_t = \mathbf{1}_{\{\eta(t)=1\}} - \lambda$$

A particular case : $E = \{1\}$

	L^1 -theory	L^2 -theory
$F \in L^2(\mathbf{P})$ s.t. $D^+F \in L^2(\mathbf{P} \otimes \nu)$	$D^+F = DF, \mathbf{P}$ -a.s.	
$u \in \text{dom} \delta$	$\tilde{\delta}(u) = \sum_t u(\eta, t) \Delta Z_t = \delta(u)$	
$F \in \text{dom} L$ $D^+F \in L^2(\mathbf{P} \otimes \nu)$	$\tilde{L}F = -\tilde{\delta}(D^+F) = -\delta(DF) = -\sum_{n \geq 1} n J_n(f_n) = LF$	

$L^1 - L^2$ correspondence?

The general case : E countable

	L^1 -theory	L^2 -theory
$F \in L^2(\mathbf{P})$ s.t. $D^+F \in L^2(\mathbf{P} \otimes \nu)$	$D^+F = DF, \mathbf{P}$ -a.s.	
u simple	$\tilde{\delta}(u) = \sum_{(t,k)} u(\eta, (t, k)) \Delta Z_{(t,k)} \neq$	$\delta(u) = \sum_{(t,k)} u(\eta, (t, k)) \Delta R_{(t,k)}$
$F \in \text{dom}L$ $D^+F \in L^2(\mathbf{P} \otimes \nu)$	$\tilde{L}F = -\tilde{\delta}(D^+F) \neq$	$LF = -\sum_{n \geq 1} n J_n(f_n) = -\delta(DF)$
$F \in L^2(\mathbf{P})$	The operator \tilde{L} is not the generator of $(P_\tau)_{\tau \geq 0}$	$P_\tau F = \sum_{n \geq 0} e^{-n\tau} J_n(f_n)$

No chance to get a L^1/L^2 unified theory

Mehler's formula and corollaries

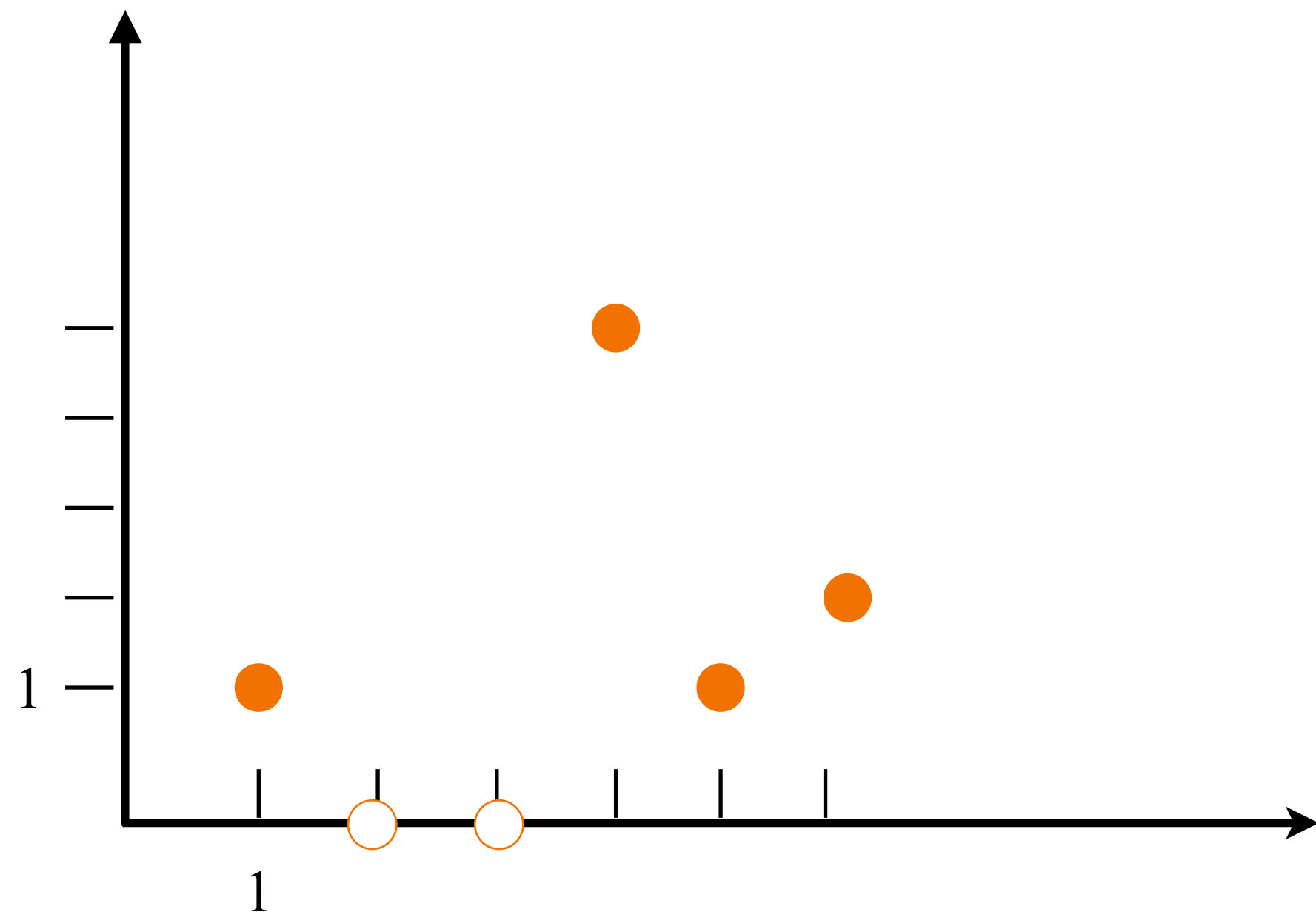
Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

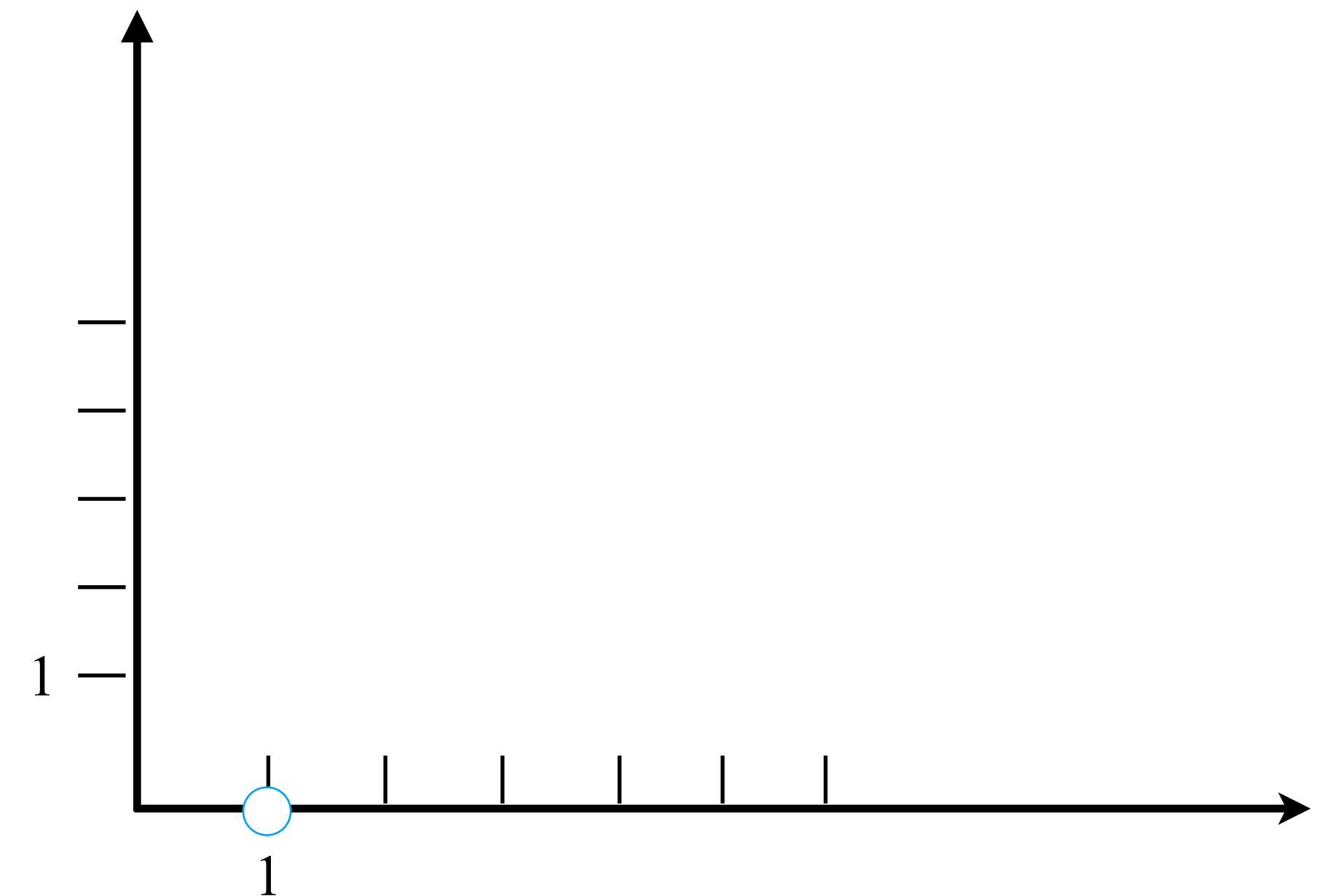
$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

Original process

$$\varepsilon_1^\tau = 1$$



Thinned process



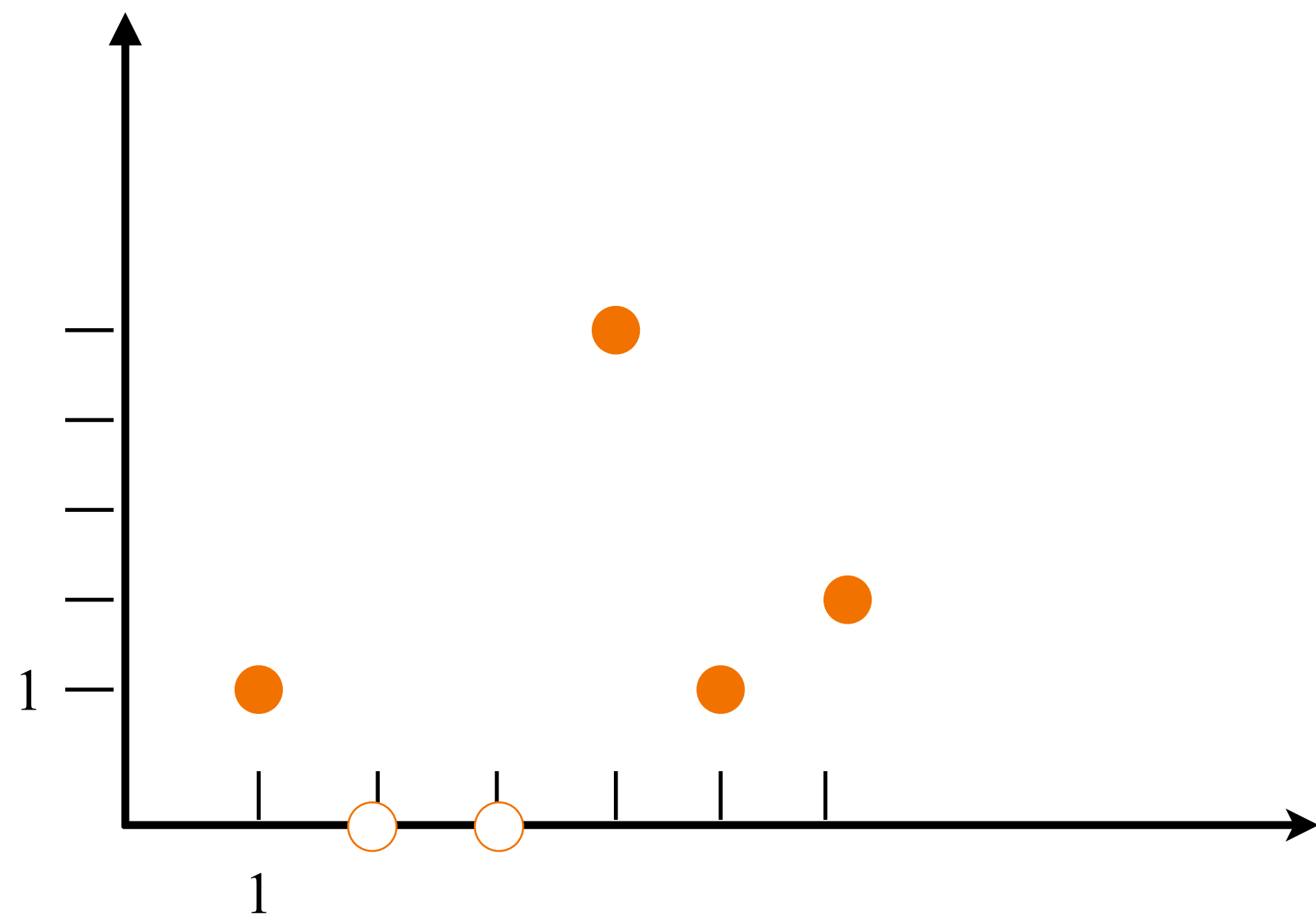
Mehler's formula and corollaries

Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

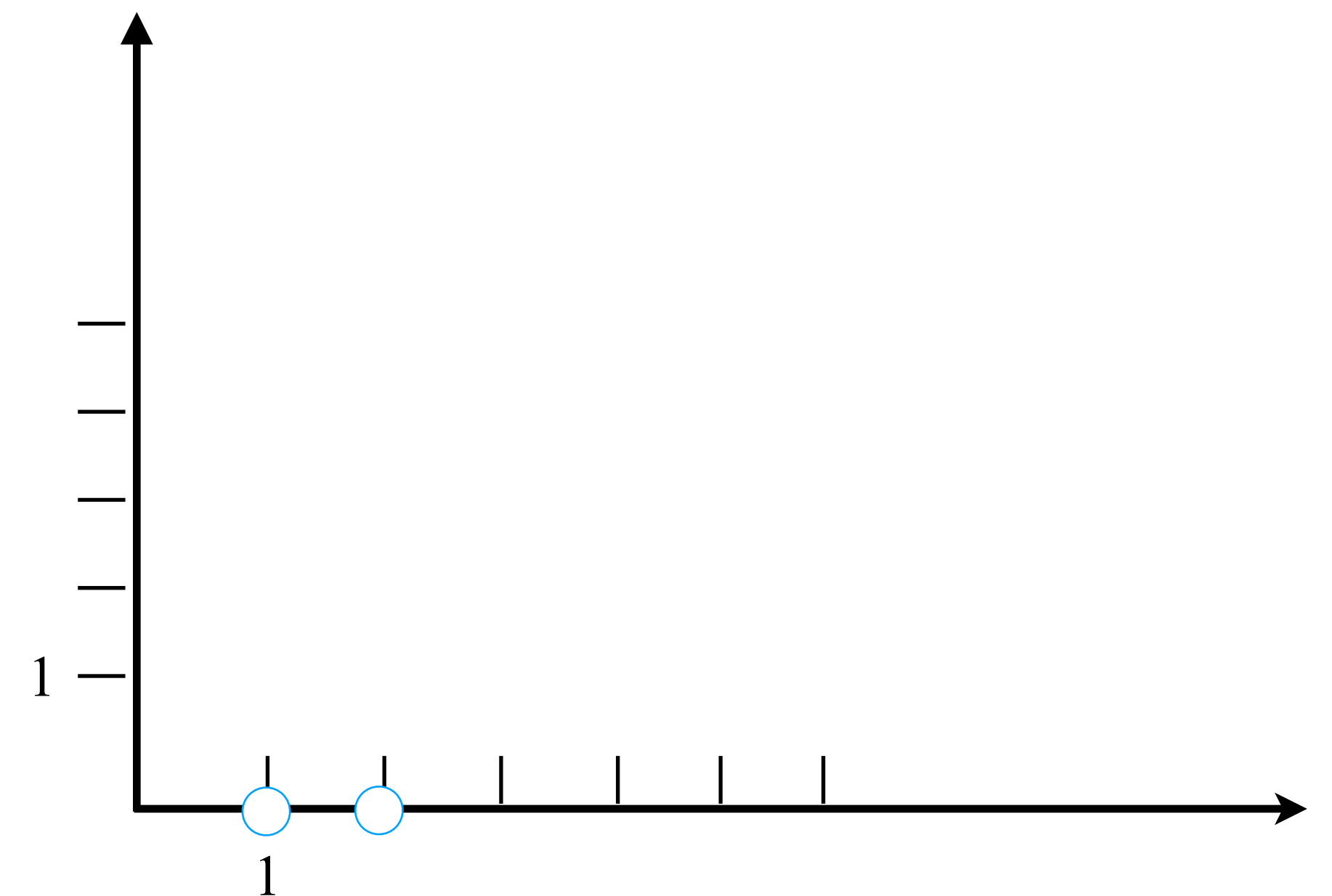
Original process



$$\varepsilon_1^\tau = 1$$

$$\varepsilon_2^\tau = 0$$

Thinned process



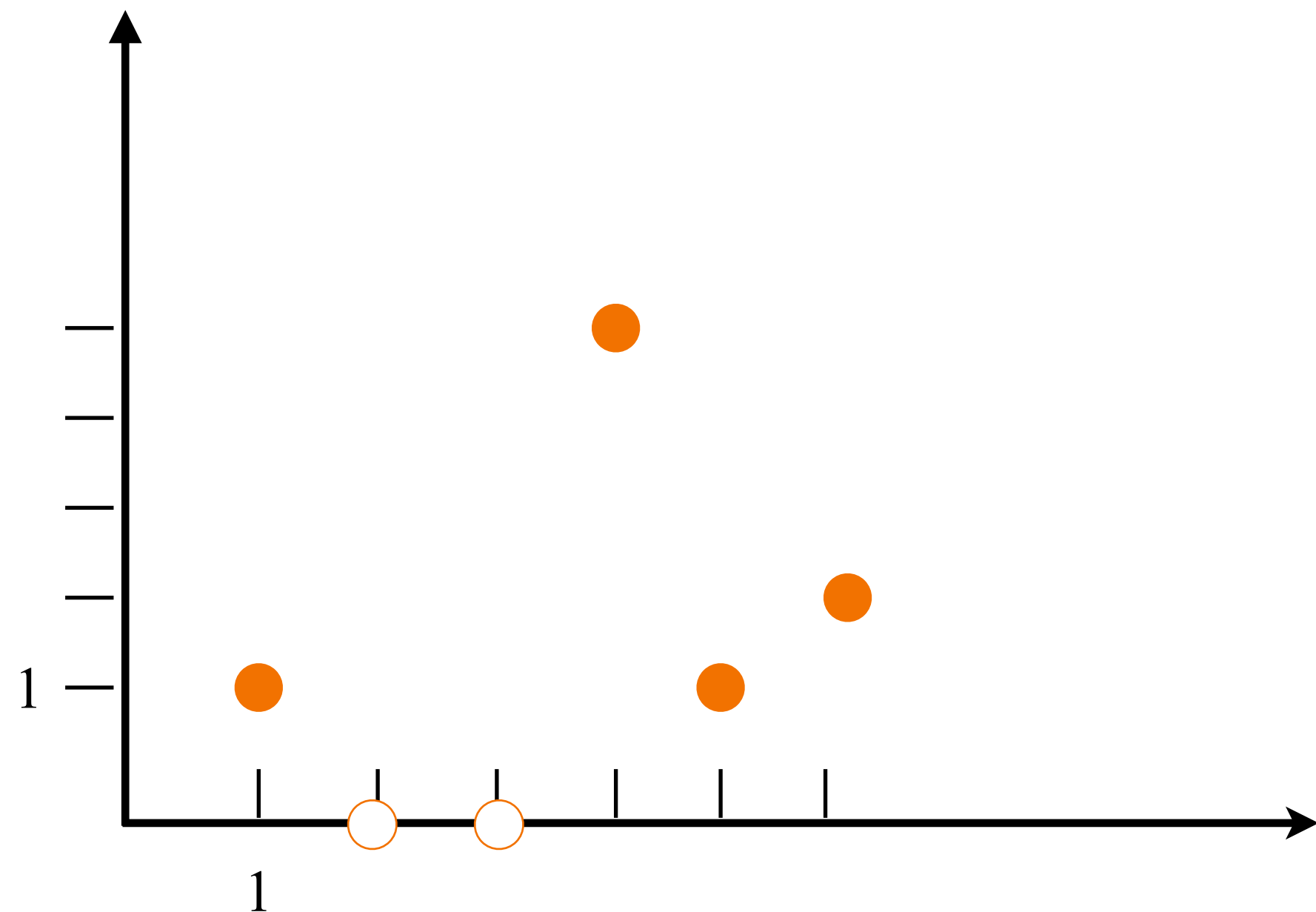
Mehler's formula and corollaries

Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

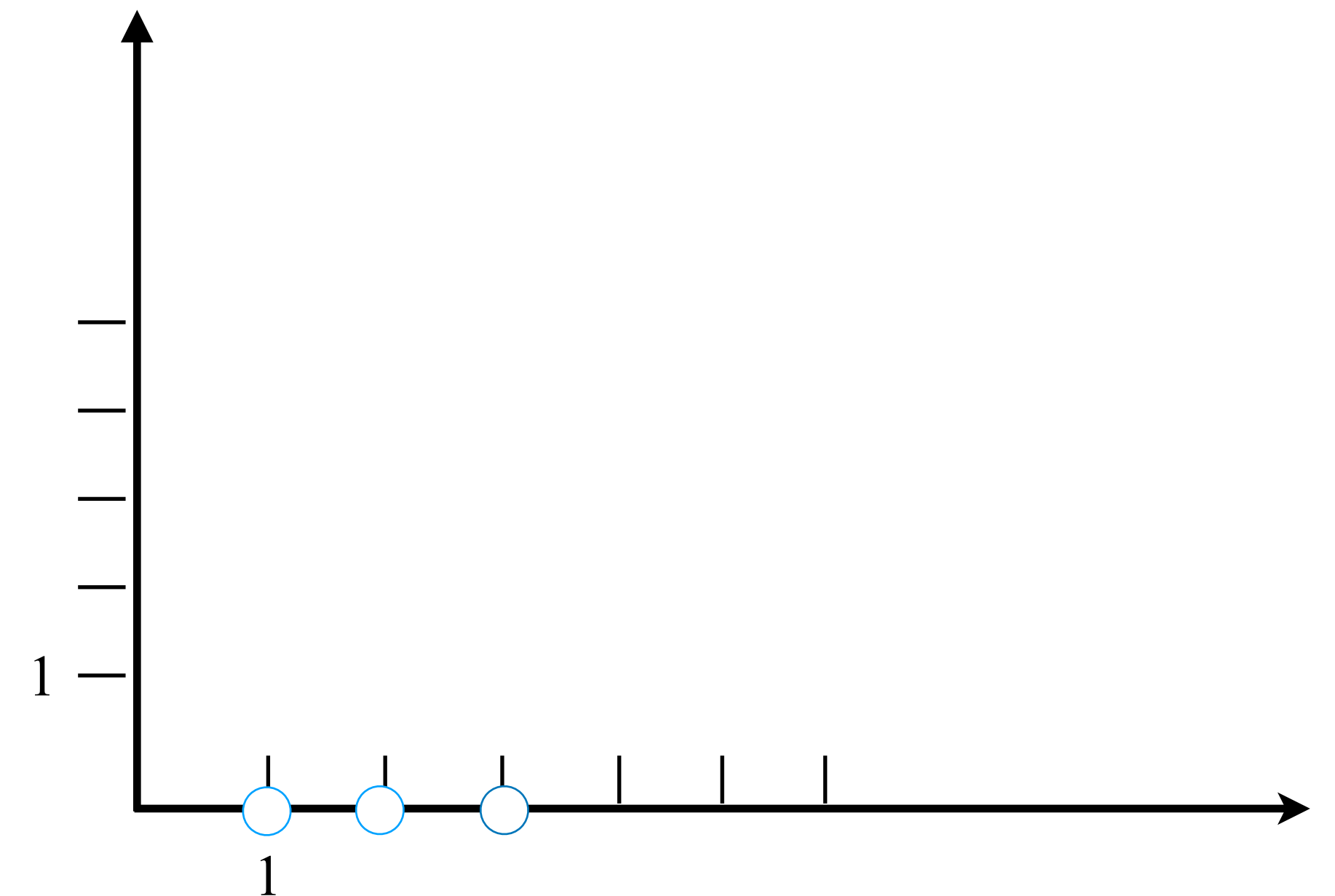
Original process



$$\varepsilon_1^\tau = \varepsilon_3^\tau = 1$$

$$\varepsilon_2^\tau = 0$$

Thinned process



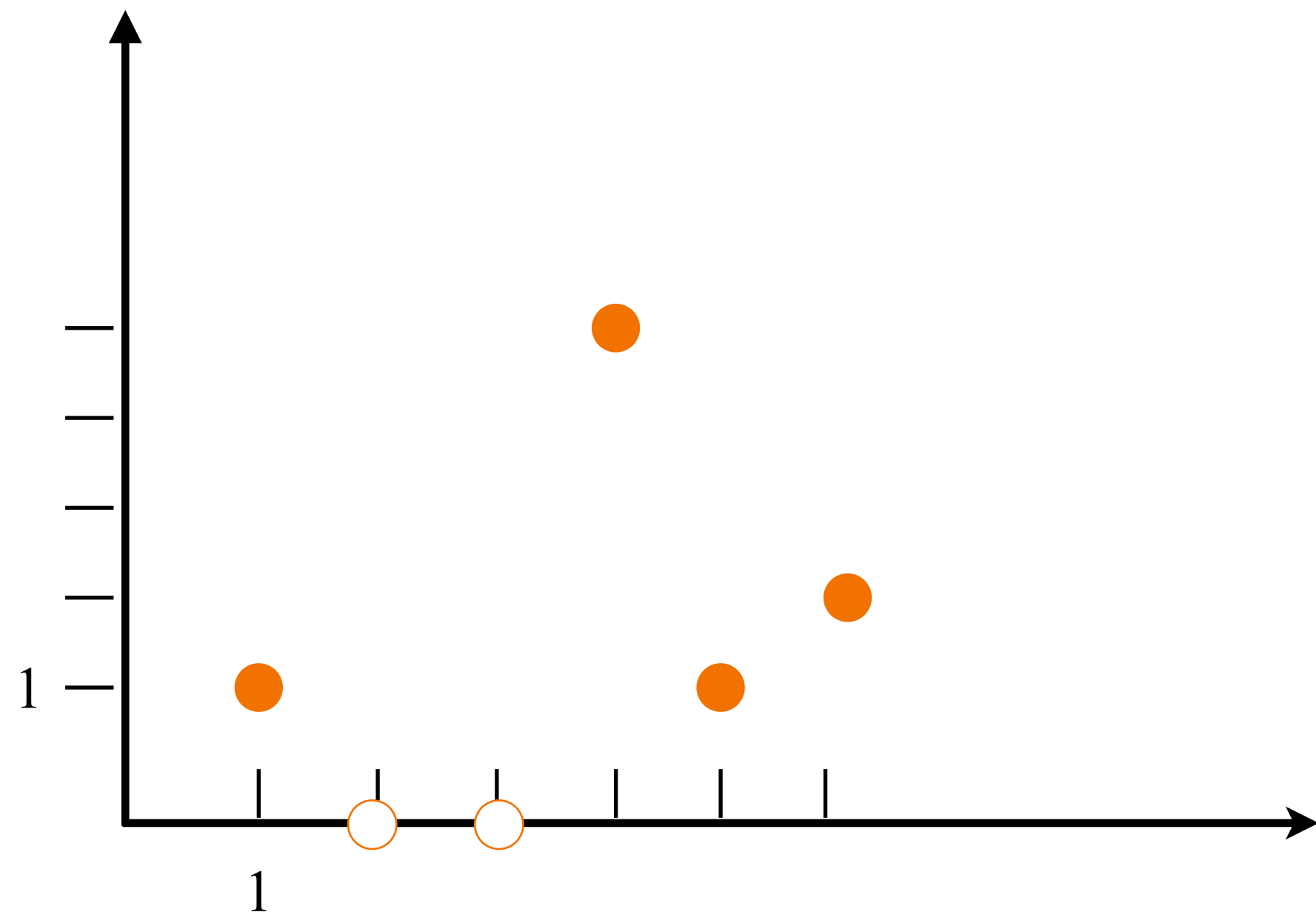
Mehler's formula and corollaries

Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

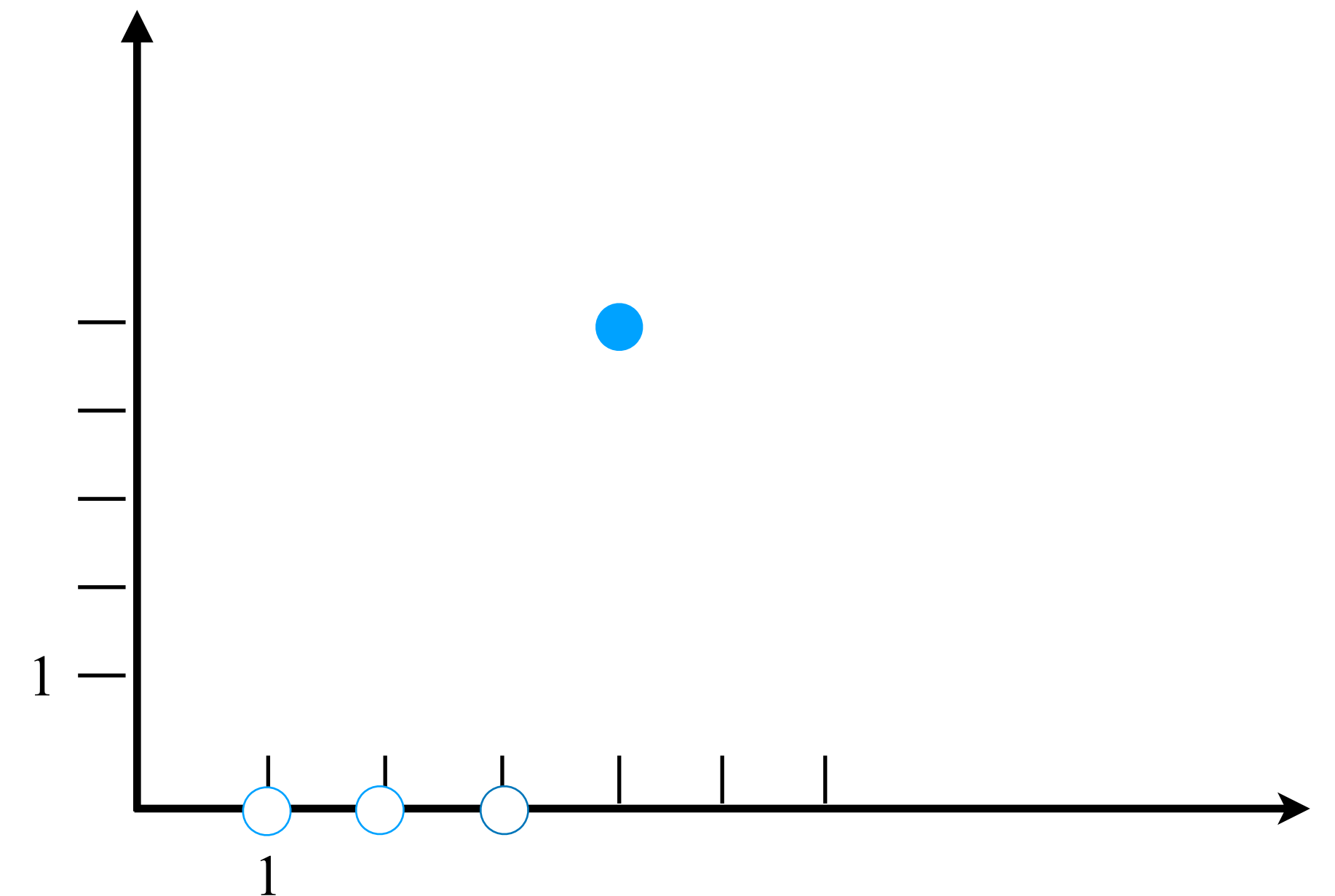
Original process



$$\varepsilon_1^\tau = \varepsilon_3^\tau = 1$$

$$\varepsilon_2^\tau = \varepsilon_4^\tau = 0$$

Thinned process



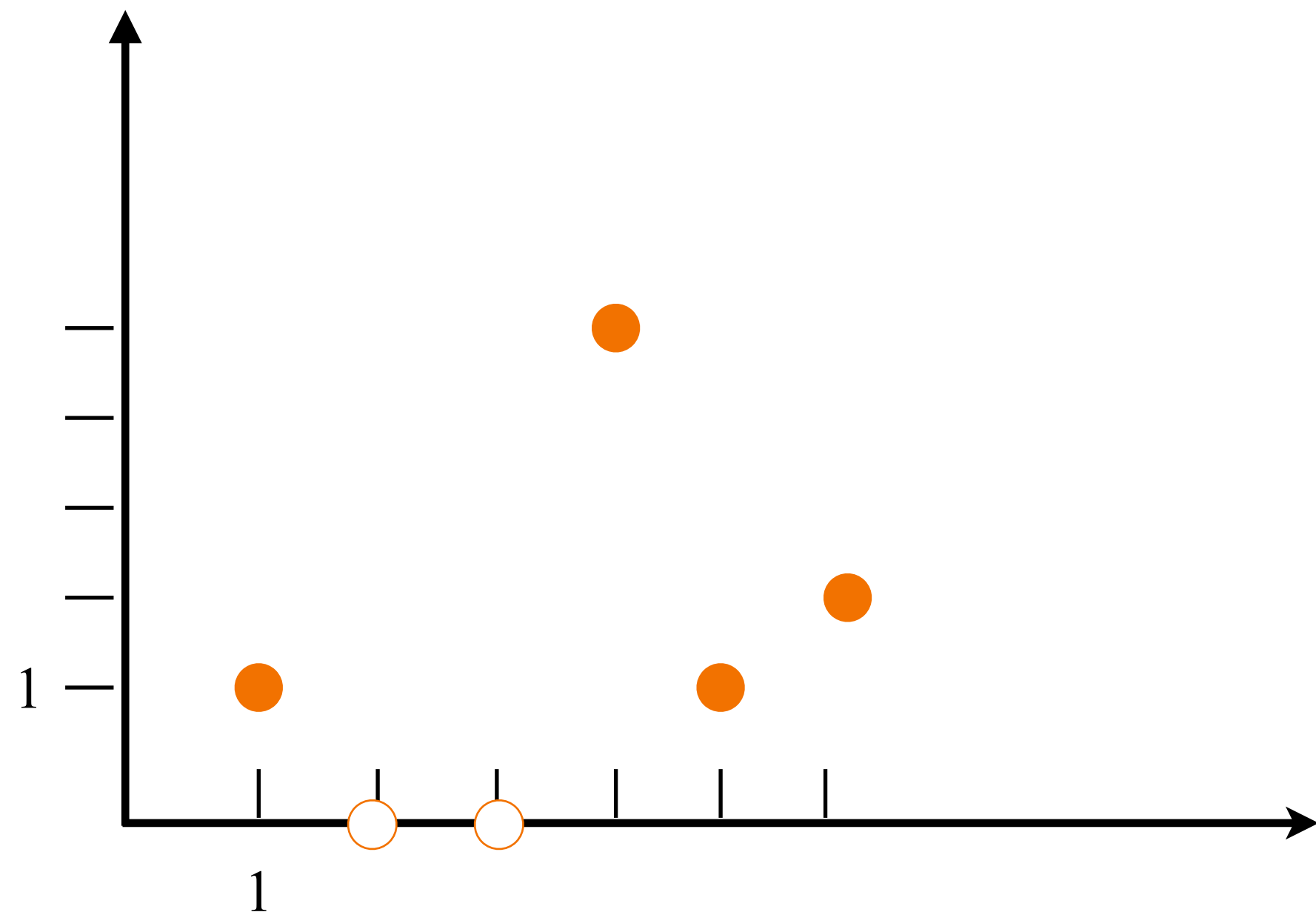
Mehler's formula and corollaries

Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

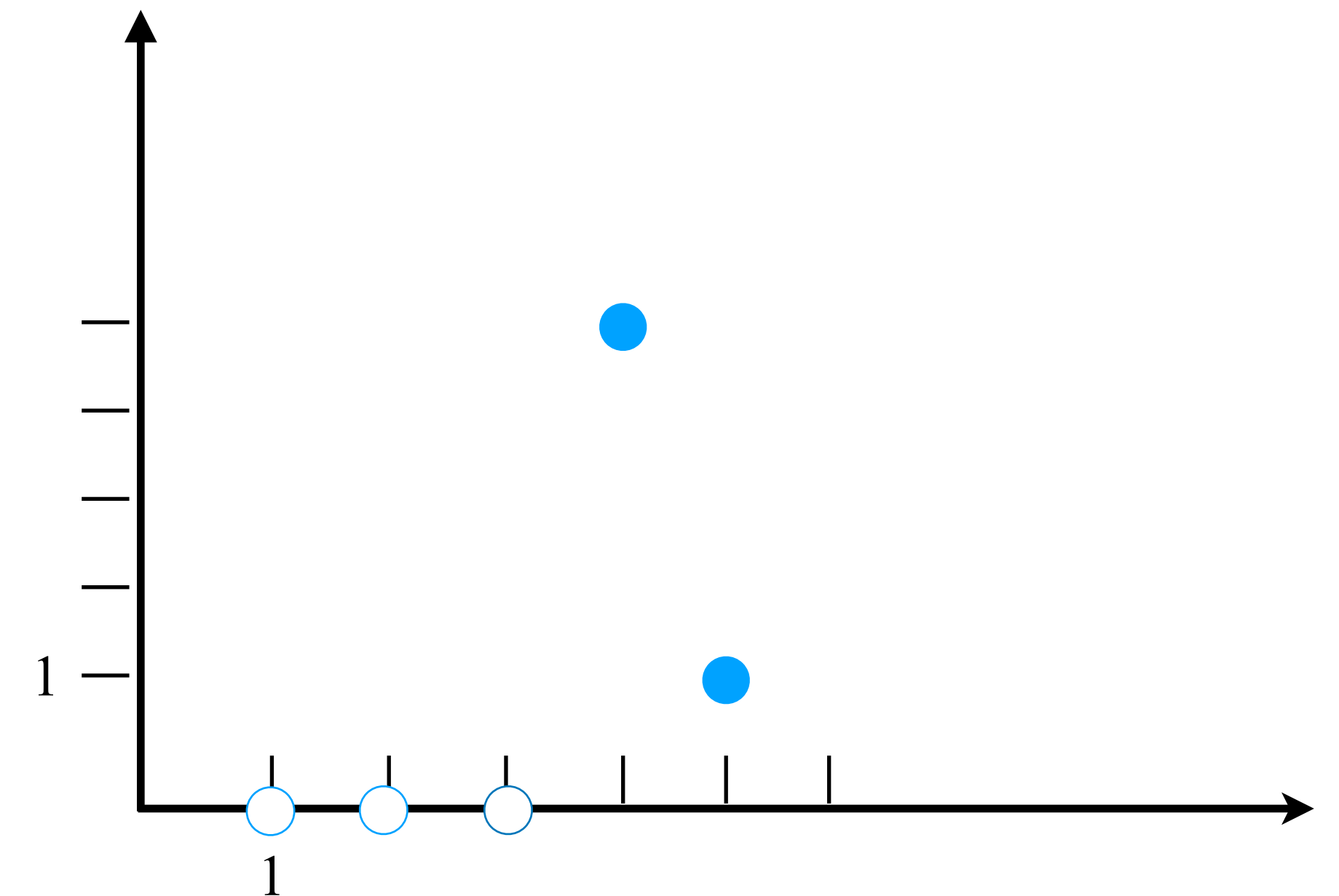
Original process



$$\varepsilon_1^\tau = \varepsilon_3^\tau = 1$$

$$\varepsilon_2^\tau = \varepsilon_4^\tau = \varepsilon_5^\tau = 0$$

Thinned process



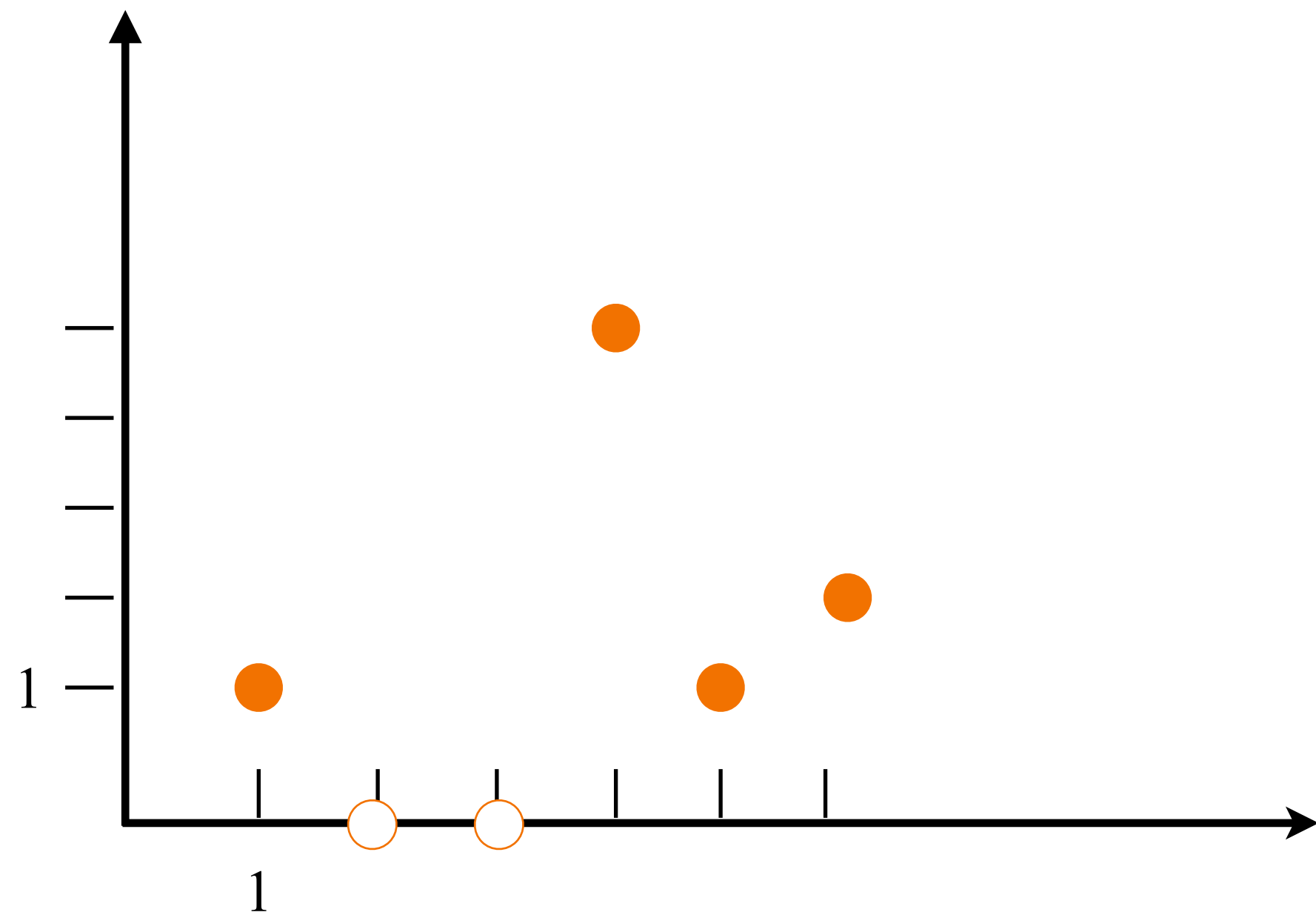
Mehler's formula and corollaries

Thinning

Let $\varepsilon_t^\tau := \mathbf{1}_{\{\theta_t \leq \tau\}}$ where $\theta_t \sim \mathcal{E}(1)$ i.i.d.

$\eta^{\tau,0} := \eta \mathbf{1}_{\{\varepsilon_t^\tau = 0\}}$ is a MBP of intensity $e^{-\tau} \nu$.

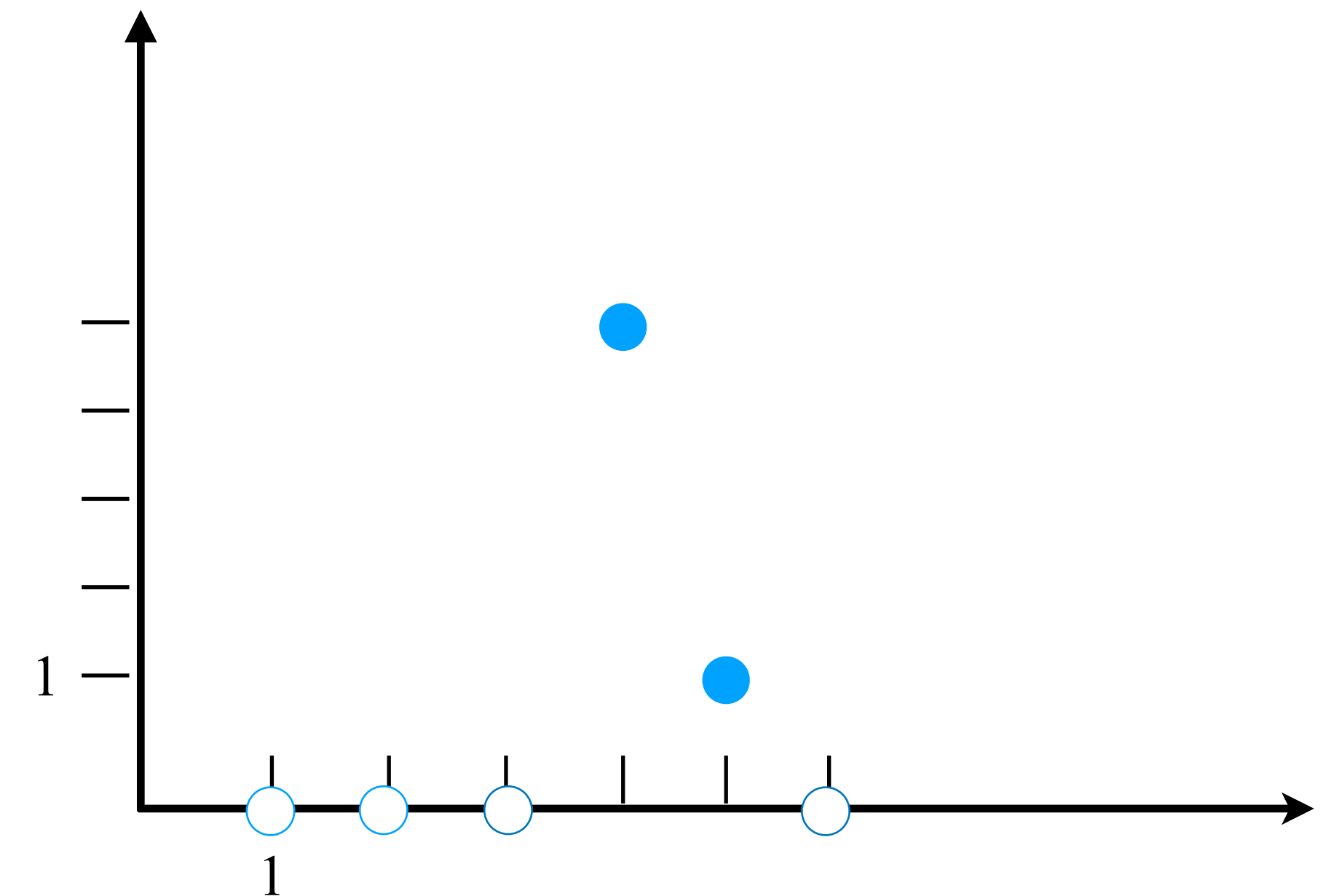
Original process



$$\varepsilon_1^\tau = \varepsilon_3^\tau = \varepsilon_4^\tau = \varepsilon_6^\tau = 1$$

$$\varepsilon_2^\tau = \varepsilon_5^\tau = \varepsilon_5^\tau = 0$$

Thinned process



Integral representation of $(P_\tau)_{\tau \geq 0}$

Proposition (H)

Let η **finite** and $F = \mathfrak{f}(\eta) \in L^1(\mathbf{P})$. For $\tau \in \mathbf{R}_+$,

$$P_\tau F = \int \mathbf{E}[\mathfrak{f}(\eta^{\tau,0} + \varepsilon \tilde{\eta}) \mid \eta] \Pi_\nu(d\tilde{\eta}) \quad \mathbf{P}\text{-a.s.}$$

where Π_ν is the distribution of a MBP(ν) and the distribution of $\tilde{\eta} \mid \eta$ is s.t.

$$\mathbf{P}((t, k) \in \tilde{\eta} \mid (t, k) \notin \eta) = \lambda \mathbf{V}(\{k\}) \quad \text{and} \quad \mathbf{P}((t, k) \notin \tilde{\eta} \mid (t, k) \in \eta) = 1 - \lambda \mathbf{V}(\{k\})$$

Integral representation of $(P_\tau)_{\tau \geq 0}$

Proposition (H)

Let η **finite** and $F = \mathfrak{f}(\eta) \in L^1(\mathbf{P})$. For $\tau \in \mathbf{R}_+$,

$$P_\tau F = \int \mathbf{E} \left[\mathfrak{f}(\eta^{\tau,0} + \varepsilon \tilde{\eta}) \mid \eta \right] \Pi_\nu(d\tilde{\eta}) \quad \mathbf{P}\text{-a.s.}$$

where Π_ν is the distribution of a MBP(ν) and the distribution of $\tilde{\eta} \mid \eta$ is s.t.

$$\mathbf{P}((t, k) \in \tilde{\eta} \mid (t, k) \notin \eta) = \lambda \mathbf{V}(\{k\}) \quad \text{and} \quad \mathbf{P}((t, k) \notin \tilde{\eta} \mid (t, k) \in \eta) = 1 - \lambda \mathbf{V}(\{k\})$$

$$\mathbf{P}((t, k) \in \tilde{\eta} \mid (t, k) \in \eta) = \mathbf{P}((t, k) \notin \tilde{\eta} \mid (t, k) \notin \eta) = 0$$

Corollaries

Corollary 1 Contraction property

For $F \in L^2(\mathbf{P})$, $\tau \in \mathbf{R}_+$,

$$DP_\tau F = e^{-\tau} P_\tau DF$$

Corollary 2 Melher's formula

For $F \in L^2(\mathbf{P})$ s.t. $\mathbf{E}[F] = 0$,

$$L^{-1}F = - \int_0^{+\infty} P_\tau F d\tau, \quad \mathbf{P} \otimes \nu\text{-a.s.} \quad \text{and} \quad DL^{-1}F = - \int_0^{+\infty} e^{-\tau} P_\tau DF d\tau,$$

Application : Poisson approximations

Poisson approximation

General bound

$$\begin{aligned}\widetilde{D}_t F &= \mathfrak{f}(\pi_t(\eta) + \delta_t) - \mathfrak{f}(\eta) \\ D_t F = D_t^+ F &= \mathfrak{f}(\pi_t(\eta) + \delta_t) - \mathfrak{f}(\pi_t(\eta))\end{aligned}$$

Theorem (H)

Let a \mathbb{N} -valued RV $F \in L^2(\mathbf{P})$ s.t. $\mathbf{E}[F] = \lambda_0 > 0$. Then

$$\begin{aligned}d_{\text{TV}}(\mathbf{P}_F, \mathcal{P}(\lambda_0)) &\leq \frac{1 - e^{-\lambda_0}}{\lambda_0} \mathbf{E} \left[\left| \lambda_0 - \langle \widetilde{D} F, -DL^{-1}(F - \mathbf{E}[F]) \rangle_{L^2(\nu)} \right| \right] \\ &\quad + \frac{1 - e^{-\lambda_0}}{\lambda_0} \mathbf{E} \left[\int_{\mathbb{N}^*} |\widetilde{D}_t F (\widetilde{D}_t F - 1) D_t L^{-1}(F - \mathbf{E}[F])| d\nu(t) \right].\end{aligned}$$

$$\nabla \varphi_A(F) = \varphi_A(F + 1) - \varphi_A(F)$$

Sketch of Proof

1. Prove that (Mehler's formula + contraction property)

$$D^+ \varphi_A(F) D^+ L^{-1}(F - \mathbf{E}[F]) \in L^1(\mathbf{P} \otimes \nu)$$

2. IBP

$$\begin{aligned} \mathbf{E} \left[F \varphi_A(F) - \lambda_0 \varphi_A(F + 1) \right] &= \mathbf{E} \left[(\tilde{L} \tilde{L}^{-1}(F - \mathbf{E}[F])) \varphi_A(F) \right] - \lambda_0 \mathbf{E}[\nabla \varphi_A(F)] \\ &= - \mathbf{E} \left[\int_{\mathbb{N}} \tilde{D}_t(\varphi_A(F)) D_t L^{-1}(F - \mathbf{E}[F]) d\nu(t) \right] - \lambda_0 \mathbf{E}[\nabla \varphi_A(F)] \\ &= - \mathbf{E} \left[\nabla \varphi_A(F) \int_{\mathbb{N}} (\tilde{D}_t F) D_t L^{-1}(F - \mathbf{E}[F]) d\nu(t) + \text{rem}_A \right] \\ &\quad - \lambda_0 \mathbf{E}[\nabla \varphi_A(F)] \end{aligned}$$

3. Bound $\mathbf{E}[|\text{rem}_A|]$ + take $\sup_{A \in \mathbb{C}\mathbb{N}}$

Application : the longest head run

$$U_n = \prod_{i=1}^{m_n} \Delta N_i + \sum_{i=1}^n (1 - \Delta N_i) \Delta N_{i+1} \Delta N_{i+2} \dots \Delta N_{i+m_n}$$

Test length $m_n = \log_{1/p}((n-1)(1-p) + 1) + \text{cst}$

Theorem (H)

Let $\lambda_n = p^{m_n}((n-1)(1-p) + 1)$. Then,

$$d_{\text{TV}}(\mathbf{P}_{U_n}, \mathcal{P}(\lambda_n)) \leq \frac{p^{m_n}}{\lambda_n} \left[p^{m_n} [2(m_n - 1)q^2 + 2m_n q + 1] + cm_n^2 (1-p)^2 p^{m_n-1} + 2(n - m_n + 1)(1-p)p^{m_n} + o(p^{m_n}) \right].$$

Remark The CDF of the longest head run can be approximated by

$$\mathbf{P}(R_n < m_n) = \mathbf{P}(U_n = 0) = e^{-\lambda_n}.$$

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, \quad C_i \sim \text{Ber}(p), \text{ i.i.d.}$$
$$U_n = \sum_{k \in I} X_k^n$$

Application : the longest head run

$$U_n = \prod_{i=1}^{m_n} \Delta N_i + \sum_{i=1}^n (1 - \Delta N_i) \Delta N_{i+1} \Delta N_{i+2} \dots \Delta N_{i+m_n}$$

$$X_k^n = (1 - C_{k-1}) \prod_{i=k}^{k+m_n-1} C_i, \quad C_i \sim \text{Ber}(p), \text{ i.i.d.}$$

$$U_n = \sum_{k \in I} X_k^n$$

Test length $m_n = \log_{1/p}((n-1)(1-p) + 1) + \text{cst}$

Theorem (H)

Let $\lambda_n = p^{m_n}((n-1)(1-p) + 1)$. Then,

$$d_{\text{TV}}(\mathbf{P}_{U_n}, \mathcal{P}(\lambda_n)) \leq \underbrace{\frac{p^{m_n}}{\lambda_n}}_{\sim 1/n} \left[p^{m_n} [2(m_n - 1)q^2 + 2m_n q + 1] + cm_n^2 (1-p)^2 p^{m_n-1} + 2(n - m_n + 1)(1-p)p^{m_n} + o(p^{m_n}) \right].$$

Remark The CDF of the longest head run can be approximated by

$$\mathbf{P}(R_n < m_n) = \mathbf{P}(U_n = 0) = e^{-\lambda_n}.$$

Compound Poisson approximation

General bound

$$\begin{aligned}\widetilde{D}_t F &= \mathfrak{f}(\pi_t(\eta) + \delta_{t,k}) - \mathfrak{f}(\eta) \\ \widetilde{L}F &= -\widetilde{\delta}D^+F \\ D_{t,k}^+ F &= \mathfrak{f}(\pi_t(\eta) + \delta_{t,k}) - \mathfrak{f}(\pi_t(\eta))\end{aligned}$$

Theorem (H)

Let $\lambda_0 > 0$, \mathbf{V} probability distribution on \mathbb{N}^* .

Let a \mathbb{N} -valued RV $F \in L^2(\mathbf{P})$ s.t. $\mathbf{E}[F] = \lambda_0 \mathbf{E}[V_1] > 0$. Then,

$$\begin{aligned}d_{\text{TV}}(\mathbf{P}_F, \mathcal{PC}(\lambda_0, \mathbf{V})) &\leq \sup_{A \subset \mathbb{Z}_+} \left| \int_{\mathbb{X}} \left[D^+ \widetilde{L}^{-1}(\mathfrak{f}(\eta) - \mathbf{E}[\mathfrak{f}(\eta)]) \psi_A(\mathfrak{f}(\pi_t(\eta) + \delta_{(t,k)}) \right. \right. \\ &\quad \left. \left. - k \psi_A(\mathfrak{f}(\eta) + k) \right] d\nu(t, k) \right| \\ &\quad + c_{\mathcal{PC}} \left| \int_{\mathbb{X}} \left[D^+ \widetilde{L}^{-1}(\mathfrak{f}(\eta) - \mathbf{E}[\mathfrak{f}(\eta)]) - k \right] d\nu(t, k) \right|\end{aligned}$$

Application : count of a rare word in a DNA sequence

$$\mathfrak{Z}(W_n) = \sum_{j \in I} Z_j := \sum_{j \in I} \mathbf{1}_{\{X_j=w_1, \dots, X_{j+h_n-1}=w_{h_n}\}}$$

Let $F_n = \sum_{j \in I} V_j \Delta N_j$ where $V_j \sim \mathbf{V} = \mathcal{G}(1 - \alpha)$ i.i.d.

$\Delta N_j \sim \text{Ber}((1 - \alpha)\mu(W_n))$ i.i.d.

Proposition (H)

Let $\lambda_n = (n - h_n + 1)(1 - \alpha)\mu(W_n)$ and $\mathbf{V} = \mathcal{G}(1 - \alpha)$. Then,

$$\begin{aligned} d_{\text{TV}}(\mathbf{P}_{\mathfrak{Z}(W_n)}, \mathcal{P}\mathcal{C}(\lambda_n, \mathbf{V})) &\leq d_{\text{TV}}(\mathbf{P}_{\mathfrak{Z}(W_n)}, \mathbf{P}_{F_n}) + d_{\text{TV}}(\mathbf{P}_{F_n}, \mathcal{P}\mathcal{C}(\lambda_n, \mathbf{V})) \\ &\leq 2h_n\mu(W_n) + (n - h_n + 1)c_{\mathcal{P}\mathcal{C}}\mu(W_n)^2 \end{aligned}$$

~~Conclusion~~ Open questions

What about inserting dependence in the process ?

Possible applications in stochastic geometry ?

References

- R. Arratia, L. Goldstein, L. Gordon, *Poisson approximation and the Chen-Stein method*, Statistical Science, 1990.
- A. D. Barbour, L. H. Y. Chen, W. L. Loh, *Compound Poisson approximation for nonnegative random variables via Stein's method*, Annals of Probability, 1992.
- G. Last, G. Peccati, M. Schulte, *Normal approximation on Poisson spaces: Mehler's formula, second order Poincaré inequalities and stabilization*, Probability Theory and Related Fields, 2016.
- G. Last, M. Penrose, *Lecture on the Poisson Process*, Cambridge University Press, 2017.
- I. Nourdin, G. Peccati, *Normal Approximations with Malliavin calculus: from Stein's method to universality*, Cambridge University Press, 2012.
- S. Schbath, *Compound approximation of word counts in DNA sequences*, ESAIM: probability and statistics, 2016.
- H. Halconrui, *Malliavin calculus for marked binomial processes and applications*, Electronic Journal of Probability, 2022.

Merci pour votre attention !