

A crash course in  
Dirichlet processes  
Part 1

Jason Swanson  
UCF Probability Seminar  
Jan 26, 2021

# I. Background and notation

①

(See Foundations of Modern Probability, Kallenberg (1997) for further details)

## RANDOM MEASURES

$(\Omega, \mathcal{F}, P)$  : Prob. sp.

$(S, \mathcal{S})$  : m'ble sp.

$M(S)$  : set of  $\sigma$ -finite measures on  $(S, \mathcal{S})$

A random measure ought to be a random variable taking values in  $M(S)$ . Need a  $\sigma$ -alg. on  $M(S)$ .

(2)

For  $A \in \mathcal{S}, B \in \mathcal{R}$ ,  $\leftarrow$  Borel  $\sigma$ -alg. on  $\mathbb{R}$

$$C_{A,B} := \{ \nu \in M(\mathcal{S}) : \nu(A) \in B \} \subset M(\mathcal{S})$$

$\uparrow$  cylinder set

$$\mathcal{M}(\mathcal{S}) := \sigma(\{C_{A,B} : A \in \mathcal{S}, B \in \mathcal{R}\})$$

$\uparrow$  a  $\sigma$ -alg. on  $M(\mathcal{S})$

$\mathcal{M}(\mathcal{S})$  is the smallest  $\sigma$ -alg. on  $M(\mathcal{S})$  that makes the projections m'ble.

(3)

For  $A \in \mathcal{S}$ , define  $\pi_A: M(S) \rightarrow \mathbb{R}$  by  
 $\uparrow$  projection onto  $A$

$$\pi_A(\nu) = \nu(A)$$

$$\mathcal{M}(S) = \sigma(\{\pi_A: A \in \mathcal{S}\})$$


---

$$M_1(S) = \{\nu \in M(S) : \nu(S) = 1\}$$

$\uparrow$  probability measures on  $S$

$$M_1(S) = C_{S, \{1\}}$$

$\uparrow S \in \mathcal{S}$        $\nwarrow \{1\} \in \mathcal{R}$

$\therefore M_1(S) \in \mathcal{M}(S)$



(4)

$$\mathcal{M}_1(S) = \{C \cap M_1(S) : C \in \mathcal{M}(S)\}$$

↖ a  $\sigma$ -alg. on  $M_1(S)$ ;  
it is  $\mathcal{M}(S)$  restricted to  $M_1(S)$

---

A *random measure* on  $S$  is an  $M(S)$ -valued random variable, i.e. a function

$\mu: \Omega \rightarrow M(S)$  that is  $(\mathcal{F}, \mathcal{M}(S))$ -m'ble.

If  $\mu(\Omega) \subset M_1(S)$ , then  $\mu$  is a *random probability measure*.

A random meas. is a special case of a "kernel".

$(T, \mathcal{T})$ : m'ble space

A **kernel** from  $T$  to  $S$  is a m'ble function  $\mu: T \rightarrow M(S)$ .

Notation: given  $t \in T$ ,  $A \in \mathcal{S}$ ,  $\mu(t, A) := (\mu(t))(A)$ .

Expl: 
$$\underset{[0, \infty)}{\overset{\mu}{t}} \mapsto \underset{\mathcal{M}([0, \infty))}{\overset{\pi}{e^{-ts} ds}} \quad (\text{need to show this is m'ble})$$

Laplace transform: 
$$F(t) = \int_{[0, \infty)} f(s) \mu(t, ds)$$

(6)

If  $\mu$  is a kernel from  $T$  to  $S$ , then  $\mu$  is a **probability kernel** if  $\mu(T) \leq M_1(S)$ .

A random (probability) measure is a (probability) kernel from  $\Omega$  to  $S$ .

Checking measurability:

$\mu: T \rightarrow M(S)$  is  $(\mathcal{I}, M(S))$ -m'ble

iff  $\pi_A \circ \mu$  is  $(\mathcal{I}, \mathbb{R})$ -m'ble  $\forall A \in \mathcal{S}$

iff  $t \mapsto \mu(t, A)$  is  $(\mathcal{I}, \mathbb{R})$ -m'ble  $\forall A \in \mathcal{S}$

## Alternative/equivalent formulation:

A **kernel** from  $T$  to  $S$  is a function  $\mu: T \times \mathcal{S} \rightarrow [0, \infty]$  such that

- $\mu(t, \cdot)$  is a  $\sigma$ -finite meas.  $\forall t \in T$ , and
  - $\mu(\cdot, A)$  is m'ble  $\forall A \in \mathcal{S}$ .
- 

## REGULAR CONDITIONAL DISTRIBUTIONS

$X$ : an  $(S, \mathcal{S})$ -valued random variable

$\mathcal{L}(X)$ : the distribution (or law) of  $X$

↑ a prob. meas. on  $S$

$$\mu = \mathcal{L}(X) \Rightarrow$$

$$\bullet P(X \in A) = \mu(A) \quad \forall A \in \mathcal{S}$$

$$\bullet E[f(X)] = \int_{\mathcal{S}} f(x) \mu(dx) \quad \text{whenever} \\ f(X) \geq 0 \text{ a.s.} \\ \text{or } E|f(X)| < \infty$$

Can we do the same thing with conditioning?

Is  $P(X \in A | \mathcal{H})$  a random prob. meas?

Can we get  $E[f(X) | \mathcal{H}]$  by integrating?

Typically, yes. But need hypotheses.

A m'ble space  $(S, \mathcal{S})$  is a (standard) Borel space if  $\exists$  a bijection  $\varphi: S \rightarrow \mathbb{R}$  such that  $\varphi$  is  $(\mathcal{S}, \mathcal{B})$ -m'ble and  $\varphi^{-1}$  is  $(\mathcal{B}, \mathcal{S})$ -m'ble.

A m'ble subset of a complete, separable metric space is a standard Borel space.

key hypothesis

$(S, \mathcal{S})$  : Borel sp.

$X$  :  $S$ -valued random var.

$(T, \mathcal{T})$  : m'ble sp.

$Y$  :  $T$ -valued random var.

Theorem  $\exists$  a kernel <sup>(probability)</sup>  $\mu$  from  $T$  to  $S$  such that

$$P(X \in A | Y) = \mu(Y, A) \text{ a.s. } \forall A \in \mathcal{S}$$

the regular conditional distribution  
of  $X$  given  $Y$ .

Notation:  $\mathcal{L}(X|Y) := \mu(Y)$

$$\mu(Y) = \mathcal{L}(X|Y) \Rightarrow$$

- $P(X \in A | Y) = \mu(Y, A)$  a.s.  $\forall A \in \mathcal{S}$
- $E[f(X) | Y] = \int_{\mathcal{S}} f(x) \mu(Y, dx)$  a.s.

whenever  $E|f(X)| < \infty$ .

What about  $X|Y$ ?

Surprisingly,  $X|Y$  is a special case of  $X|Y$ .



$\mathcal{G}$ : sub  $\sigma$ -alg. of  $\mathcal{F}$

$$(\mathcal{T}, \mathcal{G}) := (\Omega, \mathcal{G})$$

$Y$ : identity function

$$P(X \in A | \mathcal{G}) = P(X \in A | Y)$$

So  $\exists$  a prob. kernel  $\mu$  from  $\Omega$  to  $S$   
(i.e. a random prob. meas.) such that

$$P(X \in A | \mathcal{G}) = \mu(A) \text{ a.s. } \forall A \in \mathcal{S}$$

$\uparrow$  regular conditional distribution  
of  $X$  given  $\mathcal{G}$

Notation:  $\mathcal{L}(X | \mathcal{G}) := \mu$

$$\mu = \mathcal{L}(X|Y) \Rightarrow$$

- $P(X \in A | Y) = \mu(A)$  a.s.  $\forall A \in \mathcal{S}$
- $E[f(X) | Y] = \int_{\mathcal{S}} f(x) \mu(dx)$  a.s.

whenever  $E|f(X)| < \infty$ .

Helpful result:

If  $Y \in \mathcal{G}$ , then

$$E[f(X, Y) | \mathcal{G}] = \int_S f(x, Y) \mu(dx),$$

where  $\mu = \mathcal{L}(X | \mathcal{G})$ .

- Treat  $Y$  like a constant (since it is known), then integrate w.r.t. the (regular) conditional distribution of  $X$  given  $\mathcal{G}$ .
- When  $X$  &  $\mathcal{G}$  are indep.,  $\mathcal{L}(X | \mathcal{G}) = \mathcal{L}(X)$  and this result is familiar to undergrads.

Final bit of notation:

If  $\mu, \nu$  are complex measures  
(e.g. finite signed measures), then

$\mu \propto \nu$  means  $\exists c > 0$  such that

$$\mu = c\nu$$

Expl:  $X_1, X_2, \dots$  iid  $\text{Exp}(\lambda)$

$$d\mathcal{L}(X_j) \propto e^{-\lambda x} dx$$

$$d\mathcal{L}(X_1 + \dots + X_n) \propto x^{n-1} e^{-\lambda x} dx$$

notationally simpler  
to omit normalizing  
constant

## II. Laplace's sunrise problem

Take a pressed penny from a museum.

Flip it repeatedly.

$X_n$ : result of  $n^{\text{th}}$  flip (0 = tails, 1 = heads)

Assumptions:

$S = \{0, 1\}$ ,  $\mathcal{S} = \mathcal{P}(S)$  (power set)

$X_n$  is an  $S$ -valued random var.

$X_1, X_2, \dots$  are exchangeable:

$$(X_1, \dots, X_n) \stackrel{d}{=} (X_{\sigma(1)}, \dots, X_{\sigma(n)})$$

$\forall$  permutations  $\sigma$  of  $\{1, \dots, n\}$

By *de Finetti's theorem* (see Thm. 9.16 in Kallenberg):

$\exists$  a random prob. meas.  $\mu$  on  $S = \{0, 1\}$   
 such that

$\uparrow$  de Finetti says nothing  
 about the distribution  
 of  $\mu$

$$\mathcal{L}(X_1, \dots, X_n | \mu) = \mu^n \quad \forall n$$

More concisely:

$X := (X_1, X_2, \dots)$  (an  $(S^\infty, S^\infty)$ -val. r.v.)

$$\mathcal{L}(X | \mu) = \mu^\infty$$

$\mathcal{L}(X|\mu) = \mu^\infty$  means:

- $X_1, X_2, \dots$  are conditionally i.i.d. given  $\mu$ .

- $\mathcal{L}(X_n|\mu) = \mu \quad \forall n$

i.e.  $P(X_n \in A | \mu) = \mu(A)$  a.s.  $\forall n, \forall A$

$\mu$  is the (unknown) distribution of each  $X_n$ .

If we knew this distribution, flips would be i.i.d.

Without knowing  $\mu$ , flips are dependent as we learn from flip to flip.

To complete the model we must choose a distribution for  $\mu$ .

This is our "prior" on the unknown  $\mu$ , based on whatever we may know about the pressed penny before flipping it.

$\mu$  is an  $M_1(S)$ -val. r.v.

$$\mathcal{L}(\mu) \in M_1(M_1(S))$$

We must choose an element of  $M_1(M_1(S))$  to be our prior.



SIMPLIFYING  $\mu$ 

$\mu$  is an  $M_1(S)$ -val. r.v.

$$S = \{0, 1\}, \quad \mathcal{S} = \mathcal{P}(S)$$

Define  $\varphi: M_1(S) \rightarrow [0, 1]$  by

$$\varphi(\nu) = \nu(\{1\})$$

- $\varphi$  is bijective
- $\varphi$  is  $(M_1(S), \mathcal{B}_{[0,1]})$ -m'ble

← Borel  $\sigma$ -alg.  
on  $[0, 1]$

$\varphi = \pi_{\{1\}}$  is a projection

- $\varphi^{-1}$  is  $(\mathcal{B}_{[0,1]}, M_1(S))$ -m'ble

To check, use  $M_1(S) = \sigma(\text{cylinder sets})$

Define  $\theta = \varphi(\mu) = \mu(\{1\})$

$\theta$  is a  $[0,1]$ -val. r.v.

$$\sigma(\theta) = \sigma(\mu)$$

$$P(X_1=x_1, \dots, X_n=x_n | \theta) = P(X_1=x_1, \dots, X_n=x_n | \mu)$$

$$= \mu(x_1) \mu(x_2) \cdots \mu(x_n)$$

$$= \theta^a (1-\theta)^b$$

$$a = |\{j: x_j = 1\}|$$

$$b = |\{j: x_j = 0\}|$$

$\theta$  is the unknown probability of heads.

We must choose a prior distribution for  $\theta$ .

$$\mathcal{L}(\mu) \in M_1(M_1(\{0,1\}))$$

$$\mathcal{L}(\theta) \in M_1([0,1])$$

We must choose a prob. meas. on  $[0,1]$  as our prior on  $\theta$ .

Any will do, but suppose we choose one with a density.

$$d\mathcal{L}(\theta) = f(t) dt$$

posterior distribution of  $\theta$   
given first  $n$   
observations

$$\mathcal{L}(\theta | x_1, \dots, x_n) = ?$$

Basic calculations give

$$d\mathcal{L}(\theta | X_1, \dots, X_n) \propto \underbrace{t^N (1-t)^M}_{\text{beta distribution}} f(t) dt,$$

$$N = |\{j : X_j = 1\}| = \sum_{j=1}^n X_j$$

$$M = |\{j : X_j = 0\}| = n - N$$

The beta distribution is a special case of the Dirichlet distribution.

The Dirichlet distribution is a discrete version of the Dirichlet process.

# BETA DISTRIBUTION

The **beta distribution** with parameters  $\alpha > 0$  and  $\beta > 0$ , denoted **Beta( $\alpha, \beta$ )**, is the probability measure on  $[0, 1]$  proportional to

$$t^{\alpha-1} (1-t)^{\beta-1} dt.$$

$\alpha=1, \beta=1$  is the uniform distribution

If  $\mu = \text{Beta}(\alpha, \beta)$ , then

$$d\mu = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} t^{\alpha-1} (1-t)^{\beta-1} dt$$

$$= \frac{1}{B(\alpha, \beta)} t^{\alpha-1} (1-t)^{\beta-1} dt$$

↖ beta function

mean of  
Beta( $\alpha, \beta$ ) is  
 $\frac{\alpha}{\alpha+\beta}$

$$d\mathcal{L}(\theta) = f(t) dt \Rightarrow$$

$$d\mathcal{L}(\theta | X_1, \dots, X_n) \propto t^N (1-t)^M f(t) dt,$$

$$N = \sum_{j=1}^n X_j, \quad M = n - N$$

If the prior is beta, then the posterior is beta.

$$\mathcal{L}(\theta) = \text{Beta}(\alpha, \beta) \Rightarrow$$

$$\begin{aligned} d\mathcal{L}(\theta | X_1, \dots, X_n) &\propto t^N (1-t)^M t^{\alpha-1} (1-t)^{\beta-1} dt \\ &= t^{\alpha+N-1} (1-t)^{\beta+M-1} dt \end{aligned}$$

$$\Rightarrow \mathcal{L}(\theta | X_1, \dots, X_n) = \text{Beta}(\alpha+N, \beta+M)$$

Expl Flip the pressed penny  $n$  times.

Suppose we get  $s$  heads. What is the probability of heads on the  $(n+1)^{\text{th}}$  flip? Take  $\mathcal{L}(\theta)$  to be uniform?

---

$$P(X_{n+1}=1 \mid N=s) = ?$$

$$\begin{aligned} P(X_{n+1}=1 \mid N) &= E[P(X_{n+1}=1 \mid \theta, N) \mid N] \\ &= E[P(X_{n+1}=1 \mid \theta) \mid N] \\ &= E[\theta \mid N] \end{aligned}$$

$$= E[E[\theta \mid X_1, \dots, X_n] \mid N]$$

$$\mathcal{L}(\theta) = \text{Beta}(1, 1) \Rightarrow$$

$$\mathcal{L}(\theta | X_1, \dots, X_n) = \text{Beta}(1+N, 1+M)$$

$$\text{mean of Beta}(\alpha, \beta) = \frac{\alpha}{\alpha + \beta}$$

$$\therefore E[\theta | X_1, \dots, X_n] = \frac{N+1}{n+2} \quad \leftarrow N+M=n$$

$$\begin{aligned} \therefore P(X_{n+1}=1 | N) &= E[E[\theta | X_1, \dots, X_n] | N] \\ &= \frac{N+1}{n+2} \end{aligned}$$

$$P(X_{n+1}=1 | N=s) = \frac{s+1}{n+2} \quad \square$$