Quick review of some useful discrete stuff:

$\beta$ - distribution

for $x \in [0,1]$

$$P_\beta(x \mid \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot x^{(\alpha-1)} (1-x)^{(\beta-1)}$$

In this case:

$$E[X] = \frac{\alpha}{\alpha+\beta} \qquad (\text{worth remembering})$$

$\beta$ is useful for binomial problems because it has nice conjugacy properties.

Eg: I have a coin, $P(\text{heads}) = h$ (Unknown)

I toss, see $k$ heads, $n-k$ tails.

What is $P(h \mid k, n-k)$ ?

Natural choice: $\quad P(h) = \beta(\alpha_\pi, \beta_\pi)$

$\uparrow$

prior

then $\qquad$ $P(h \mid k, n-k) \propto P(K, n-k \mid h) \, P(h)$

$$\begin{array}{c} \uparrow \\ \text{binomial} \\ \left\{ \cdot h^{K} (1-h)^{n-k} \right. \end{array}$$

$$\begin{array}{c} \uparrow \\ \beta \\ h^{\alpha_{\pi}-1} \cdot (1-h)^{\beta_{\pi}-1} \end{array}$$

So product

$$\propto \qquad h^{(\alpha_{\pi}+K-1)} (1-h)^{(n-k+\beta_{\pi}-1)}$$

which is $\qquad \beta\left(\alpha_{\pi}+K, \ \beta_{\pi}+n-k\right)$

This is one of those cute things one learns and forgets, but there is a more general point

## Dirichlet

assume we have
$$x_{1} \cdots x_{K}$$

on a <u>simplex</u>

i.e. $\quad 0 < x_{i} < 1$

$$x_{1} + \cdots x_{K} = \underline{1}$$

$\qquad$ and $\quad \alpha_{1} \cdots \alpha_{K} > 0$

Then

$$P(x_1 \cdots x_K \mid \alpha_1 \cdots \alpha_K) = \frac{1}{B(\alpha)} \prod_i x_i^{(\alpha_i - 1)}$$

Dirichlet

(Notice how this generalizes $\beta$ dist - its a PDF on Prob distributions )

Again, there are neat conjugacy properties

Eg: Assume I roll a K-sided die N times, observing faces $N_1 \cdots N_K$ times.

~~cont~~ What is $P(P_1 \cdots P_K \mid n_1 \cdots n_K)$?

Bayesian: Prior on $P_1 \cdots P_K$ is $Dir(\alpha_1 \cdots \alpha_K)$

then $P(P_1 \cdots P_K \mid n_1 \cdots n_K) \propto P(n_1 \cdots n_K \mid P_1 \cdots P_K) P(P_1 \cdots P_K)$

multinomial

$\{ P_1^{n_1} P_2^{n_2} \cdots P_K^{n_K}$

$P_1^{\alpha_1 - 1} P_2^{\alpha_2 - 1} \cdots$

This gives

$$P(P_1 \cdots P_k \mid n_1 \cdots n_K) \text{ is } Dir(\alpha_1 + n_1, \cdots \alpha_K + n_*)$$

Sometimes written as

$$\alpha = (\alpha_1 \cdots \alpha_K) \qquad \text{concentration hyperparameter.}$$

prob a
face comes
up $\rightarrow$  $p \mid \alpha \qquad \sim \qquad Dir(K, \alpha)$

$X \mid p \qquad \sim \qquad Cat(K, p)$

$\underset{\text{multinomial}}{\curvearrowleft}$

$\uparrow$
observed
faces

then $\quad c = (c_1 \cdots c_K) = $ counts of cat $1 \cdots K$

$$p \mid X, \alpha \quad \sim \quad Dir(K, \alpha_1 + c_1 \cdots \alpha_K + c_K).$$

Now, one fair objection to previous methods is that we have no probabilities

— We can **bound** expected counts below, but not compute expectations.

Consider

$$p \sim Dir(\alpha_1, \cdots \alpha_K)$$

$$E(p) = \frac{(\alpha_1, \cdots \alpha_K)}{\sum_{j=1}^{K} \alpha_j}$$

(which is a discrete probability distribution, as it should be)

Now assume we do not know anything about which face of die is more likely

$\longrightarrow$ suggests a prior

$$Dir(\alpha) \quad \text{where} \quad \alpha = \alpha_0 \cdot 1$$

(i.e. all $\alpha$'s the same).

In each case, mean will be uniform pd.

but

$$\text{Dir}(\alpha) \quad \propto \quad x_1^{\alpha-1} \ldots x_K^{\alpha-1}$$

So for very large $\alpha$, few $x_i$'s can be large (they can't all be smaller than $\frac{1}{K}$!)

— <u>concentration parameter</u>

$\alpha$    big    $\Rightarrow$    concentrated dist's

     small    $\Rightarrow$    more diffuse.

To compute expectations, we need a random model of $CdG(\lambda)$.

Desirable object:

· a way of constructing random probability distributions, so that we can still compute expectations, etc.

There are several such objects

· Dirichlet process:

$X$ — a Polish space

= separable, completely metrizable topological space.

(examples:

Real line, Interval, etc)

= space homeomorphic to a complete metric space w/ a ~~some~~ countable dense subset;

$\alpha$ . a finite measure.

i.e. $\alpha(X) < \infty$.

$P$ — a <u>random</u> measure.

i.e. "randomly chosen" and a measure.

then $P$ is a Dirichlet process
if for every finite measurable
partition $\{B_1 \cdots B_K\}$. of $X$,

the Joint of $\{P(B_1), \cdots P(B_K)\}$ is

a $K$-dimensional Dirichlet dist with
params.

$$\alpha(B_1), \quad \cdots \quad \alpha(B_K)$$

$\alpha$ is called the <u>base measure</u>

Notice the property is quite strong; isn't obvious one exists.

Notice also that "big" sets in $\alpha$ ~~get are~~

tend to be "big" in $P$.

$$\boxed{\text{Work on } \mathbb{R}}$$

Some properties:

• consider a partition $\{A, A^c \equiv \mathbb{R} - A\}$

$$P(A) \text{ is } \beta(\alpha(A), \alpha(A^c)).$$

$\underset{\text{random variable}}{\uparrow} \qquad \underset{\text{distribution}}{\uparrow}$

So $E[P(A)] = \dfrac{\alpha(A)}{\alpha(A) + \alpha(A^c)}$

Now write $\alpha = M \cdot G$

$$\underset{\substack{\mathbb{R} \\ \alpha(\mathbb{R})}}{\big\uparrow} \qquad \overset{\text{prob measure}}{\nwarrow}$$

$\therefore E[P(A)] = G(A)$

Now assume $P$ is a D.P.

and

$$X|P \sim P$$

└─────────────┘
↑ Yona can think of $P$ as
a prior $P(X|P) = P$

then the <u>marginal</u> of $X$ is

$G_k$  ⟵ this is a <u>distribution</u>

[Show this from the expectations]

Now recall $P(A) \sim \beta(\alpha(A), \alpha(A^c))$.

Variance of $\beta(\alpha, \beta) = \dfrac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

So

$$Var(P(A)) = \frac{\alpha(A)\alpha(A^c)}{M^2(M+1)} = \frac{G(A)G(A^c)}{(M+1)}$$

$\Rightarrow$ bigger $M$ implies "more concentration".

Notice that easy arguments give

$$E\left[\int \psi \, dP\right] = \int \psi \, dG$$

from above. This is useful, because it allows us to think about integrals, as above.

P is distributed as D.P. with base measure $\alpha = MG$, write $P \sim DP(M, G)$

## Conjugacy :

assume $X_1 \cdots X_n \mid P \sim P$, IID

$$P \sim DP(M, G) = D_\alpha$$

then

$$P \mid X_1 \cdots X_n \sim D_{\alpha + \sum_{i=1}^{n} \delta(X_i)}$$

Interpret

$$D_{\alpha + \sum_i \delta(x_i)} \qquad \text{as} \qquad \text{another DP,}$$

with base measure $\qquad \alpha + \sum_i \delta(X_i)$

$\qquad\qquad\qquad\qquad\qquad\qquad \llcorner \delta$ functions
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ at $X_i$

This must be true, because:
(weaker than proof!)

Consider partition $(A_1, \cdots A_r.)$

and counts $\qquad n_1 \cdots n_r$

$$n_j = \sum_i 1(X_i \in A_j)$$

$P(A_1) \cdots P(A_r)$ is Dirichlet; prior is $\cancel{\alpha(A)}$

$$Dir(r, \alpha(A_1) \cdots \alpha(A_r))$$

$\therefore$ posterior must be

$$Dir(r, \alpha(A_1) + n_1 \cdots \alpha(A_r) + n_r)$$

And this must work for all partitions

In turn,

$$E(P \mid X_1 \cdots X_n) = \frac{M}{M+n} G + \frac{n}{M+n} \cdot \mathbb{P}_n$$

empirical dist,
is at each data point.

It can be proven that:

ⓐ      samples from the DP are discrete measures, with prob 1

ⓑ Dirichlet processes have nice self similarity properties

let $A$ be some set st
$$0 < G(A) < 1$$

$P|_A$  is  the restriction of $P$ to $A$

i.e.  $P|_A (B) = \dfrac{P(A \cap B)}{P(A)}$

Now  $P|_A$  is  $DP(MG(A), G|_A)$

$P|_{A^c}$  is  $DP(MG(A^c), G|_{A^c})$

(You can see' this by thinking about marginals.)

and they are independent.

AND

$P|_A$  is independent of  $P(A)$

so this means that, for any given
set A.

- how mass is dist. in A
- .  .  .  .  in $A^c$
- the total mass of A

are independent.

Predictive distributions from DP.

consider $\qquad P \sim DP$

$$X_1 \cdots X_n \sim P \qquad IID$$

1) $\qquad X_1 \sim G$ $\qquad$ ( ~~by exp~~ marginal,

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ from expectation

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ result. )

2) $X_2 \mid \cancel{P}, X_1 \sim P$

$\qquad$ and $\qquad P \mid X_1 \sim DP\left(M+1, \dfrac{M}{M+1} G + \dfrac{1}{M+1} \delta X_1\right)$

$\qquad\qquad\qquad\qquad\qquad$ (from result above on DP and

$\qquad\qquad\qquad\qquad\qquad\qquad$ data)

Now Take expectation OVER $P$ :

$$X_2 \mid X_1 \sim \dfrac{M}{M+1} G + \dfrac{1}{M+1} \delta X_1$$

i.e.

- with prob $\frac{1}{M+1}$, $X_2$ is duplicate of $X_1$

- otherwise, it's New.

We can extend this argument

$$X_n | X_1 \cdots X_{n-1}$$

(we will have $n_i$ copies of $\theta_i$ where $\theta_i$ is one of the <u>distinct</u> values)

So
$$X_n | X_1, \cdots X_{n-1} \sim \begin{cases} \delta\theta_i & \text{prob } \dfrac{n_i}{M+n-1} \\[4mm] G & \text{prob } \dfrac{M}{M+n-1} \end{cases}$$

Now $p(\text{new value})$ at Step 1 is 1

| | | |
|---|---|---|
| | 2 | $\dfrac{M}{M+1}$ |
| | 3 | $\dfrac{M}{M+2}$ |
| | $n$ | $\dfrac{M}{M+n-1}$ |

so write $K_n$ for number of <u>distinct</u> values

$$E(K_n) = \sum_{i=1}^{n} \frac{M}{M+i-1}$$

This looks like $M \log \frac{n}{M}$ as $n \to \infty$

(i.e. few new things.

    — <u>great as a clustering prior</u>.

    — not like words, objects     )

<u>Stick breaking repn.</u>

Consider      $\theta_i \sim G$    IID.

              $Y_i \sim \beta(1, M)$    IID

Write

$$P = \sum_{i=1}^{\infty} V_i \, \delta\theta_i$$

where    $V_1 = Y_1$ ,   $V_2 = (1-Y_1) Y_2$ ,   $V_3 = (1-Y_1)(1-Y_2)Y_3$

$$V_i = Y_i \prod_{j=1}^{i-1} (1-Y_j)$$

Stick breaking because $Y_1$ breaks off
part of stick, $Y_2$ breaks off part of
remains $[(1-Y_1)]$, etc.

notice this has a <u>recursive</u> form

$$P \stackrel{d}{=} Y_1 \delta_{\theta_1} + (1-Y_1) P$$

$\uparrow$ Sample from $DP(M, G)$    $\uparrow$ Sample from $DP(M, G)$

This makes it quite easy to prove some
things

1) $E\left(\int \psi \, dP\right) = \int \psi \, dG$

$$\int \psi \, dP = Y_1 \cdot \psi(\theta_1) + (1-Y_1) \int \psi \, dP$$

so $E\left(\int \psi \, dP\right) = E(Y_1 \psi(\theta_1)) + E\left((1-Y_1) \int \psi \, dP\right)$

But: $Y_1, \theta_1$ are <u>indep</u>

$(1-Y_1), \int \psi \, dP$ are <u>indep</u>.

So

$$E(\int \Psi dP) = E(Y_1) \, E\left(\Psi_t(\theta_1)\right)$$

$$+ \quad E(1-Y_1) \cdot E(\int \Psi dP).$$

$$Y_1 \sim \beta(1, M) \qquad So \qquad E(Y_1) = \frac{1}{M+1}$$

$$E(1-Y_1) = \frac{M}{M+1}$$

So

$$\left(1 - \frac{M}{M+1}\right) E(\int \Psi dP) = \frac{1}{M+1} \, E(\Psi_1(\theta_1))$$

$$\uparrow$$

$$\int \Psi dG \quad \text{by weak}$$

$$\text{law of large}$$

So $\quad E(\int \Psi dP) = \int \Psi dG \quad$ nums.

a similar argument yields $\quad \text{Var}(\int \Psi dP).$

Now assume we have

$$G \sim DP(\alpha, H) = DP_{\alpha H}$$

$\alpha H$ is base measure $\quad H(A) = 1 \quad$ is prob. dist.

we know $G | X_1 \cdots X_n \sim DP\left(\alpha + n, \dfrac{\alpha H + \sum_i \delta X_i}{\alpha + n}\right)$

from earlier arguments.

<u>Another representation</u>:

consider $\quad G | X_1 \cdots X_n . \qquad = P$

~~$P(X_i)$~~ there are $K$ unique $X_i$, $K < n$

call these $X_1^* \cdots X_K^*$

then $P(X_1^*), P(X_2^*), \quad \cdots \quad P(X_K^*), P\left(A - \sum_i X_i^*\right)$

$\sim Dir\left(n_1, \quad \cdots \quad n_K, \left[\dfrac{\alpha H\left(A - \sum_i X_i^*\right)}{\alpha + n}\right]\alpha + n\right)$

$\xcancel{=} Dir\left(n_1, \quad \cdots \quad n_K, \alpha\right).$

This means we can
represent

$$G|_{x_1, \cdots x_n} \quad \text{as}$$

$$\sum_{i=1}^{K} P_i \, \delta X_i^* \quad + \quad P_{K+1} \, G^*$$

where
$$\left( P_1 \quad \cdots \quad P_{K+1} \right) \sim Dir(n_1 \cdots n_K, \alpha)$$

and $\quad G' \sim DP(\alpha, H)$.

a DP won't really do the work for us.

~~2 parameter~~

• Richer family of random processes

    • 2 parameter Poisson-Dirichlet process

    OR.    Pitman-Yor process

Choose:

    Base PDF $G$     —

$$0 < \sigma < 1$$
$$\theta > -\sigma$$
$$\Big] - parameters$$

Now choose $x_i \overset{iid}{\sim} G$

$$V_1 \sim \beta(1-\sigma, \theta+\sigma)$$
$$V_i \sim \beta(1-\sigma, \theta+i\sigma)$$

$$P = \sum_{i=1}^{\infty} \tilde{p} \, \delta(x_i) \qquad \text{where } \hat{p} = V_1$$

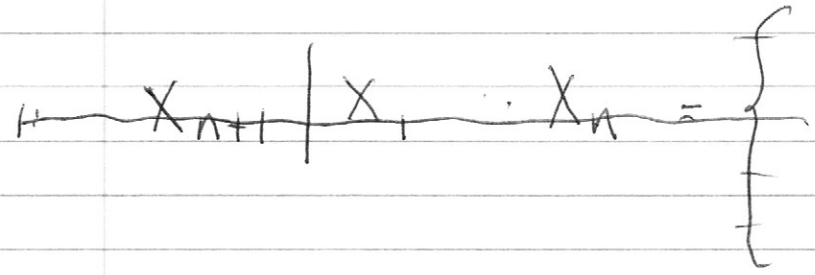$$\tilde{p}_i = V_i \prod_{i=1}^{i-1}(1-V_j) \qquad i \geq 2$$

This object $P$ is a sample from a Pitman-Yor process.

Notice $\tilde{p}_i$ is a series

1) $\tilde{p}_i \geqslant 0$, $\sum_{i=1}^{\infty} \tilde{p}_i = 1$ a.s

Proving facts about $P$ involves manipulating PDF's on such sequences (which is hard).

- State results, but do not prove them.

$$\vdash\!\!-\!\!-\!\!X_{n+1}\Big|X_1 \cdots X_n = \Bigg\{$$

- assume we see $X_1 \cdots X_n$
- notice we should have duplicates
- write $X^*_{n_{ij}}$ for $j^{th}$ distinct of $k$ species
- $n_{n_{ij}}$ for # of this type

Then

$$X_{n+1} \mid X_1 \cdots X_n = \begin{cases} x_i & \text{with prob } \dfrac{(\theta + \sigma k_n)}{(\theta + n)} \quad \swarrow \text{ a new one} \\[4mm] X^*_{n,j} & \text{"} \quad \dfrac{(n_{n,j} - \sigma)}{(\theta + n)} \end{cases}$$

◦ Notice

- rich get richer

 i.e. if you've seen many $X^*_{n,j}$,
 you'll see more.

- prob of seeing a new species
 is high if $k_n$ similar to $n$

(Not sure how to prove this — proof is
associated w/ Pitman + Yor )

Notice also this gives us a coverage result.

Coverage

$$= \% \text{ of species rep.}$$

$$= 1 - (\text{prob next one will be new})$$

$$= 1 - \left( \frac{\theta + \sigma k_n}{\theta + n} \right)$$

this could be quite high

eg. Shakespeare data.

$$\frac{\theta + \sigma(37k)}{\theta + \sigma(885k)}$$

$\theta = 1$  $\sim 0.05$  | $\theta$ matters!

$\theta = 10^6$  $\sim 0.5$

Now consider drawing n values from

$$P \sim PD(\sigma, \theta).$$

what is <u>posterior</u> on P ?

assume $K \le n$ distinct values

$$X_1^* \cdot \quad X_K$$

with $n_1 \cdots n_K$ ~~freq's~~. counts.

Then

$$P \mid X \stackrel{d}{=} \sum_{j=1}^{K} P_j^* \delta_{X_j^*} + \left(1 - \sum_{j=1}^{K} P_j^*\right) \tilde{P}^{(K)}$$

*there are errors here ! see 2] a*

where $P_1^*, \quad P_j^* \sim Dir\left(\dfrac{n_1 - \sigma}{\theta + K\sigma}, \quad \cdots \dfrac{n_K - \sigma}{\theta + K\sigma}\right)$

and $\tilde{P}^{(K)}$ is $PD(\sigma, \theta + K\sigma)$

$$P \mid X \overset{d}{=} \sum_{j=1}^{k} P_j^* \, \delta_{x_j^*} + R^* \, \tilde{P}^{(K)}$$

$$\text{where} \quad P_1^* \cdots P_K^*, \; R^* \sim \text{Dir}\left(n_1 - \sigma, \cdots n_K - \sigma, \; \Theta + K\sigma\right)$$

$$\text{and} \quad \tilde{P}^K \text{ is } PD\left(\sigma, \; \Theta + K\sigma\right).$$

One way to view this model is as a probability distribution on partitions.

Partition of $n = \sum_{i=1}^{K_n} N_{i,n}$ things into $K_n$ partitions each containing $N_{i,n}$ things

$$( \quad )( \quad ) \cdots ( \quad )$$

$$N_{1,n} \qquad N_{2,n} \qquad\qquad N_{K_n, n}$$

Notice previous results give ~~P(new par~~ expressions for where a new element goes, given an existing partition.

From these, we can derive

$$P\left(K_n = k, N_{1n} = n_1, \cdots N_{K_n n} = n_K\right)$$

$$= \frac{\prod_{i=1}^{K-1} (\theta + i\sigma)}{(\theta+1)_{n-1}} \cdot \prod_{j=1}^{K} (1-\sigma)_{n_j - 1}$$

Using the notation

$$a_b = a \cdot (a+1) \cdots (a+b-1)$$

and.

$$a_0 = 1$$

This yields a strategy for estimating $\theta, \sigma$

$$\max \quad P\left(K_n = k, N_{1n} = n_1, \cdots N_{K_n} = n_K \mid \theta, \sigma\right)$$

( max likelihood )

$\rightarrow$ should give large $\theta$ for

many small categories