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

<u>Natural choice</u>: $P(h) = \beta(\alpha_\pi, \beta_\pi)$

$\uparrow$

prior

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

$\uparrow$

$\underset{\text{binomial}}{\phantom{x}}$

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

$\uparrow$

$\beta$

$h^{\alpha_{\pi} - 1} \cdot (1-h)^{\beta_{\pi}^{-1}}$

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. $\quad 0 < x_i < 1$

$$x_1 + \cdots x_K = \underline{1}$$

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

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

Sometimes written as

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

prob a
face comes $\quad p | \alpha \quad \sim \quad Dir(K, \alpha)$
up

$$\nearrow X | p \quad \sim \quad Cat(K, p)$$

$$\underset{\text{multinomial}}{\uparrow}$$

observed
faces

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

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

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

    — We can <u>bound</u> 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

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

Soo 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>

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

$\quad$ small $\quad \Rightarrow \quad$ more diffuse.

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

Desirable object:

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

There are several such objects

· Dirichlet process:

    $\underline{X}$ — a Polish space

        = separable, completely metrizable topological space.

(examples:

Real line, Interval, etc )

        = space homeomorphic to a complete metric space w/ a ~~same~~ countable <u>dense</u> subset:

$\alpha$ .      a finite measure.

i.e. $\alpha(\mathcal{X}) < \infty$ .

$P$ — 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 $\mathcal{X}$,
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 $\{A, A^c \equiv \cancel{\mathbb{R}}-A\}$

$$P(A) \text{ is } \beta(\alpha(A), \alpha(A^c)).$$
$$\underset{\underline{\text{random variable}}}{\uparrow} \qquad \underset{\underline{\text{distribution}}}{\uparrow}$$

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

Now write $\alpha = M \cdot G$
$$\underset{\underset{\alpha(\cancel{\mathbb{R}})}{\mathbb{R}}}{\Bigg\uparrow} \quad \overset{\text{prob measure}}{\nwarrow}$$

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

Now assume $P$ is a D.P.

and
$$X|P \sim P$$

└─ You can think of $P$ as a prior $P(X|P) = P$

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

$$G_k \qquad \leftarrow \text{ this is a } \underline{\text{distribution}}$$

[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

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

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