

①

- We want to represent counts

Poisson R.V.'s capture the idea of uniformity

- The count in an interval
 - depends on length of interval
 - in fact, \propto length
- Events are independent.

(Idea works in multiple dimensions).

If X is Poisson, intensity λ

$$\bullet E(X) = \lambda$$

$$\bullet \text{Var}(X) = \lambda$$

Notice λ has units (eg $\frac{\#}{s}$, etc)

and depends on scale of interval.

PMF :

$$P(X = n \mid \text{unit interval, } \lambda) = \frac{e^{-\lambda} \cdot \lambda^n}{n!}$$

straightforward series manipulation gives

- this is a PMF
- expectation
- variance

Now assume whatever we're watching is • Poisson

- types are independent

[These assumptions are a stretch; words aren't like this; neither are animals or objects; but they're simple + generic]

Notice:

- i) If I observe a Poisson RV for an interval longer than 1

$$P(X=n \mid \text{Interval length } t, \lambda) = e^{-\lambda t} \frac{(\lambda t)^x}{x!}$$

- 2) If I observe Poisson RV for N unit intervals, seeing count n_i in the i 'th, Max likelihood est of λ is

$$\lambda^* = \frac{1}{N} \sum_i n_i$$

- 3) If I observe Poisson RV for N ~~unit~~ intervals, i th has length L_i and see n_i in i 'th, Max likelihood gives

$$\lambda^* = \frac{1}{N} \sum_i \left(\frac{n_i}{L_i} \right)$$

③

Now assume we see $\left\{ \begin{array}{l} \text{words} \\ \text{animals} \\ \text{objects} \end{array} \right\}$ for

an interval (which could be time or space). We choose a scale so this interval is $[-1; 0]$.

Each word type has intensity

λ ← intensity
 w_0 ← type

which is unknown

Natural to try and draw conclusions from word type counts.

$n_{sc} = \left[\begin{array}{l} \text{number of word types} \\ \text{that appear } x \text{ times} \end{array} \right]$

Natural because Max likelihood on individual words is no help.

~~MF~~

3a

• Also, assume future is like the past.

ie. if a word has λ_w in $(-1, 0]$
it has this λ later.

• E+T phrase this as
"conditionally binomial"

④

• ML est. of λ for words we haven't seen is 0

• But n_1 large compared to n_2 , etc suggests there are words types where

× we haven't seen them

× $\lambda_w > 0$

× So we should be looking at $G(\lambda)$.

$$G(\lambda) = P(\lambda_w \leq \lambda)$$

↑ clearly, a discrete distribution

• represents CDF (cumulative dist. function)

$$dG(\lambda) \equiv p(\lambda_w = \lambda) d\lambda$$

↑ s functions or atoms

(6)

Now

$$P(\text{a } \uparrow \text{ given word type has count } x \mid \lambda_w)$$

$$= e^{-\lambda_w} \frac{\lambda_w^x}{x!}$$

$$P(\text{a } \uparrow \text{ given word type has count } x)$$

$$= \int e^{-\lambda} \frac{\lambda^x}{x!} dG(\lambda)$$

$$\text{or } \int e^{-\lambda} \frac{\lambda^x}{x!} p(\lambda) d\lambda$$

$$E[\# \text{ of word types w/ count } x]$$

$$= \sum P(\text{word type } i \text{ has count } x)$$

i.e. word types

$$= C \int e^{-\lambda} \frac{\lambda^x}{x!} dG(\lambda) = \eta_x$$

↑ total # of word types, unknown!

(7)

- Another way to look at this is that

$d\Gamma(\lambda) = C dG(\lambda)$ is a measure
(like a PDF, +ve, but doesn't \int to 1)

- Notice that C could be hard to get, because we could have support for $G(\lambda)$ at, say,

$$\lambda = 10^{-12}$$



- there is a word we see about once in 10^{12} intervals.
- affects C , but not a significant effect on observations.

8

Now

- η_x is the observed value of an RV
- call it r_x
- we have

- $E(r_x) = \eta_x$

- reasonable approx

?

- $\text{Var}(r_x) = \eta_x$

Sum of Poisson

$$r_x \sim N(\eta_x, \sqrt{\eta_x})$$

Sum of random variables

- This will come in useful.

(9)

Now consider

$$\mathbb{E} \Delta(t) = \mathbb{E} \left[\begin{array}{l} \# \text{ of types seen in} \\ [0; t], \\ \text{but not in } [-1; 0] \end{array} \right]$$

then

$$\Delta(t) = \int_0^{\infty} \left[\begin{array}{l} e^{-\lambda} \lambda^0 \\ \lambda \\ 0! \end{array} \right] \left[1 - e^{-\lambda t} \right] dG(\lambda)$$

seen 0 times
in $[-1, 0]$

not seen 0 times in
 $[0, t]$

Notice you can derive expressions for

$$\mathbb{E} \left[\begin{array}{l} \# \text{ seen } a \text{ times in } [-1, 0] \text{ and} \\ b \text{ times in } [0, t] \end{array} \right]$$

in the straightforward way.

- Notice also that you don't need C to evaluate this, just $CdG(\lambda)$.
- Q: estimate $\Delta(t)$ given n_{oc}

• Notice

$$1 - e^{-\lambda t} = \lambda t - \frac{(\lambda t)^2}{2!} + \frac{(\lambda t)^3}{3!} - \frac{(\lambda t)^4}{4!} \dots$$

So:

$$\Delta(t) = \eta_1 t - \eta_2 t^2 + \eta_3 t^3 - \eta_4 t^4 \dots$$

(assuming convergence, etc).

Natural estimator:

assume $\eta_1 = n_1$, $\eta_2 = n_2$, etc.
 substitute