

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.

(2)

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

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

- 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 ← 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

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

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

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

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

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

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

$$= \sum_{i \in \text{word types}} P(\text{word type } i \text{ has count } x)$$

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

↑ 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

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

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

11

Paper gives word counts, word type counts

884647 words,

14376 only once

31534 word types

Q: if we saw another 884647 words,
how many new words would there
be?

A: $t = \frac{\# \text{ of words in new}}{\# \text{ in ref}}$

$= 1$

$\Delta(t) = 11430 = \Delta(1)$

$\text{Var}(\Delta(i))$ by assuming the n_x
are indep, poisson

$\text{Var}\{\Delta(i)\} \approx \sum_{i=1}^{\infty} n_i t^i = 31534$

std = 178

(Skip Fisher model)

Notice that, for $t > 1$

$$n_1 t - n_2 t^2 + n_3 t^3 - \dots$$

is a series that oscillates savagely.

You could interpret this several ways

- it doesn't converge. — panic
- the oscillations "cancel", and we need some way to accelerate this cancellation

↑
 quite plausible, as $n_x \rightarrow 0$
 as $x \rightarrow \infty$.

this gives § 4 — Euler's transform.

Now skip to § 7.

(13)

Recall that $\Delta(t)$ may be hard to estimate for large t (because there may be low frequency words).

Instead, try for a lower bound on

$\Delta(t)$ — call this bound

$b(t)$

Problem becomes:

$$b(t) = \inf_{CdG(\lambda)} \left[\int_0^{\infty} e^{-\lambda} [1 - e^{-\lambda t}] [CdG(\lambda)] \right]$$

Subject to

$$\eta_x = \mathbb{E} \int_0^{\infty} \begin{bmatrix} e^{-\lambda} x \\ \lambda \\ x! \end{bmatrix} [CdG(\lambda)]$$

This would be an LP in $CdG(\lambda)$,

IF we knew η_x .

Strategy 1:

* assume $n_x = n_{oc}$
 * discretize n_x

→ then we have an honest LP and can solve.

* This works for $\mathbb{E} + T$, but failed on object data (infeasible - ?!?)

Strategy 2:

* assume ~~n_x~~ $n_x \leq n_{oc} + \gamma \sqrt{n_{oc}}$

$$n_x \geq n_{oc} - \gamma \sqrt{n_{oc}}$$

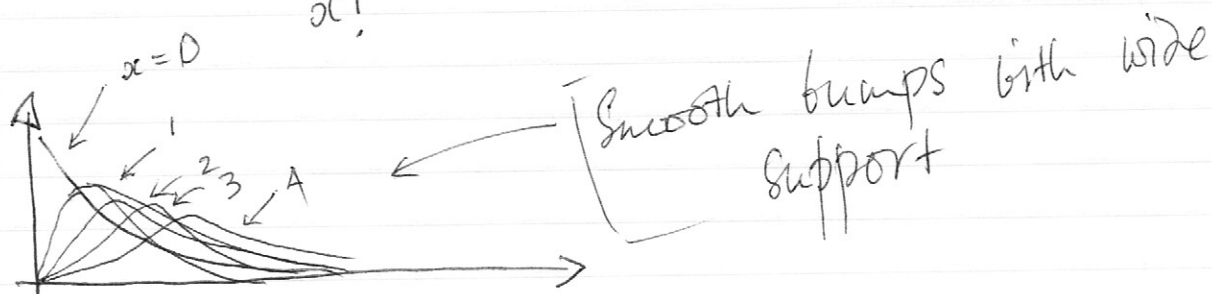
(i.e. γ stds away from mean)

* This might slacken the bound, but is probably better practice

but S_2 isn't all that reliable either
 (on objects, get feasibility ~~is~~ only if you
 apply big λ AND use only
 n_x for $x \in [1 \dots 5]$ which is worrying.)

What is going on here?

Consider $\frac{e^{-\lambda} \lambda^x}{x!}$ as a function of λ

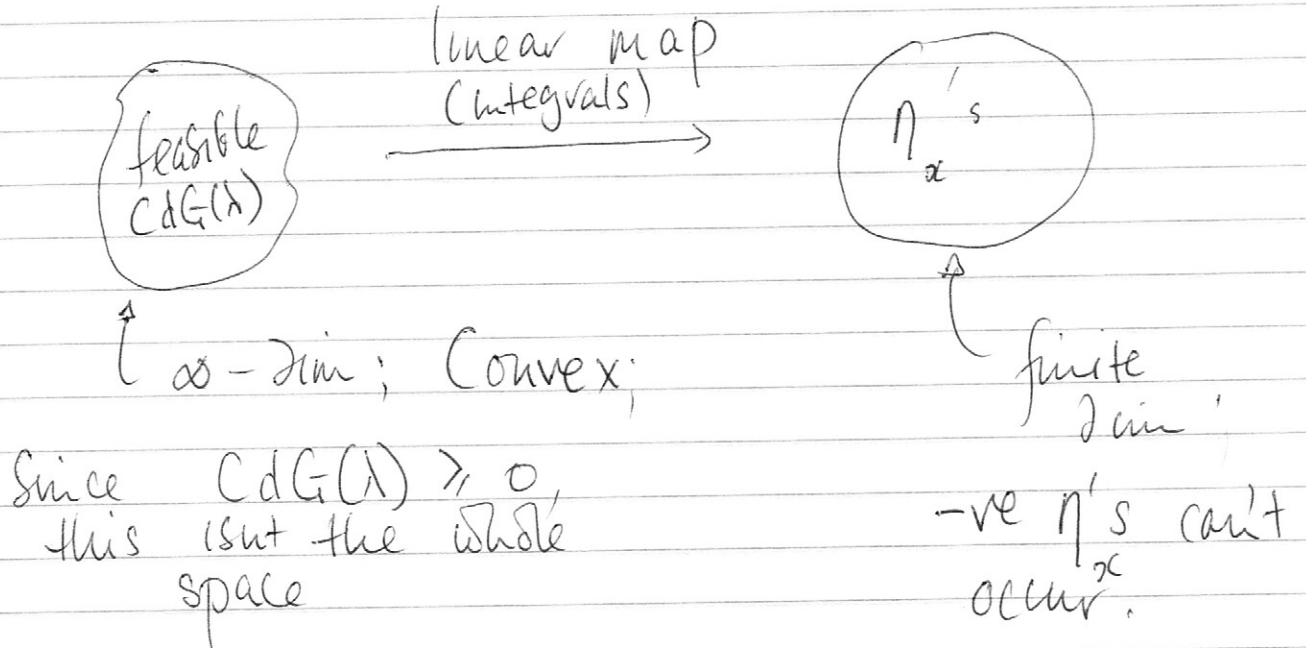


This means we can't have

$$\eta_1 = 100; \eta_2 = 0; \eta_3 = 100; \text{ etc.}$$

because these functions overlap so
 strongly $\rightarrow \int \frac{e^{-\lambda} \lambda^i}{i!} \text{CdG}(\lambda)$ is similar to
 $\left(\frac{e^{-\lambda} \lambda^i}{i!} \text{CdG}(\lambda) \right)$

Alternative View:



~~This picture strongly implies that~~

use this picture with plots;

there are vectors of η_x

that are (a) non-negative

(b) infeasible

and if you use $\eta_x = \eta_x$,

this gets infeasibility

Also explains why large r is required,
AND why large x creates problems
(the n_x estimates are poor).

What to do?

we actually know quite a lot about r_x ,
which is what we should be working

~~#~~ r_x with

- approximately Gaussian
- var approx = mean.

Strategy 3

assume ~~#~~ $r_x \sim N(n_x, \sqrt{n_x})$

and $n_x = r_x$.

Now we must (a) estimate the counts and (b) estimate $b(t)$

$$\inf_{CdG(\lambda)} \int e^{-\lambda} [1 - e^{-\lambda t}] [CdG(\lambda)] + \mu \sum_{i=1}^K S_i^2$$

↑
weight

$$S_i^2 = \frac{(r_i - n_i)^2}{2 n_i^2} \quad \leftarrow \text{Standard error of } i\text{th count,}$$

$$n_i = \int \frac{e^{-\lambda} \lambda^i}{i!} C dG(\lambda) \quad \leftarrow \text{Count estimates}$$

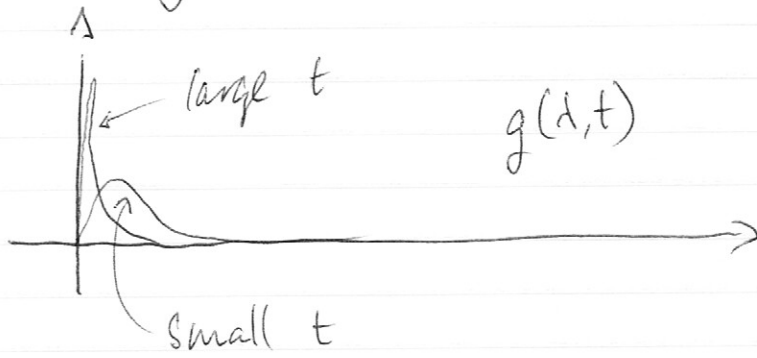
Q: how to choose μ ?

A: cross validation

Now a QP, but it's convex, so no worries

Notice also that

$$g(\lambda, t) = e^{-\lambda} (1 - e^{-\lambda t}) \text{ is important}$$



- This (basically) looks for weight in small λ 's (of $CdG(\lambda)$) which makes sense.

Issue: • lower bounds are helpful, but we want more (estimates, etc).

• Options

- work w/ continuous $CdG(\lambda)$ models

- make models explicitly discrete.

- Notice one attractive feature of this ~~problem~~ formulation
 - γ_x , $b(t)$, $\Delta(t)$ are quite insensitive to "small" changes in $Cd(GD)$
 - because
$$e^{-\lambda} \frac{\lambda^x}{x!}$$
 is pretty smooth.
- ~~This~~ This means that big shifts in where the weight is in a prob model are required to change counts
- In turn, we could use quite rough discretizations (but finer near 0).