

(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

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

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

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

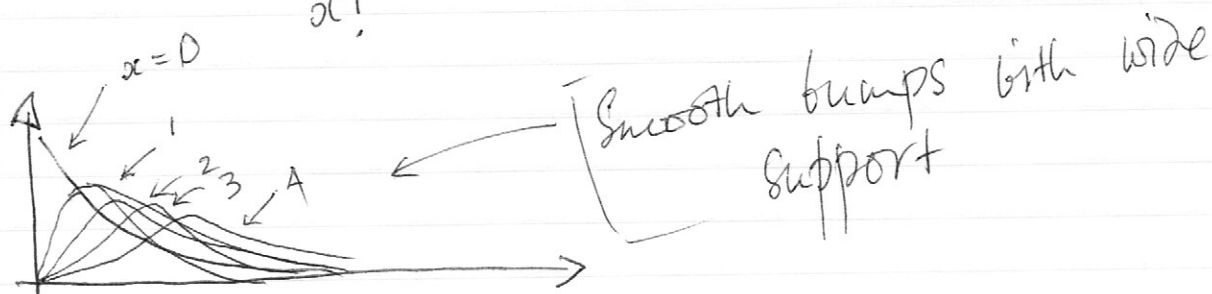
(i.e. γ stds away from mean)

* 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?

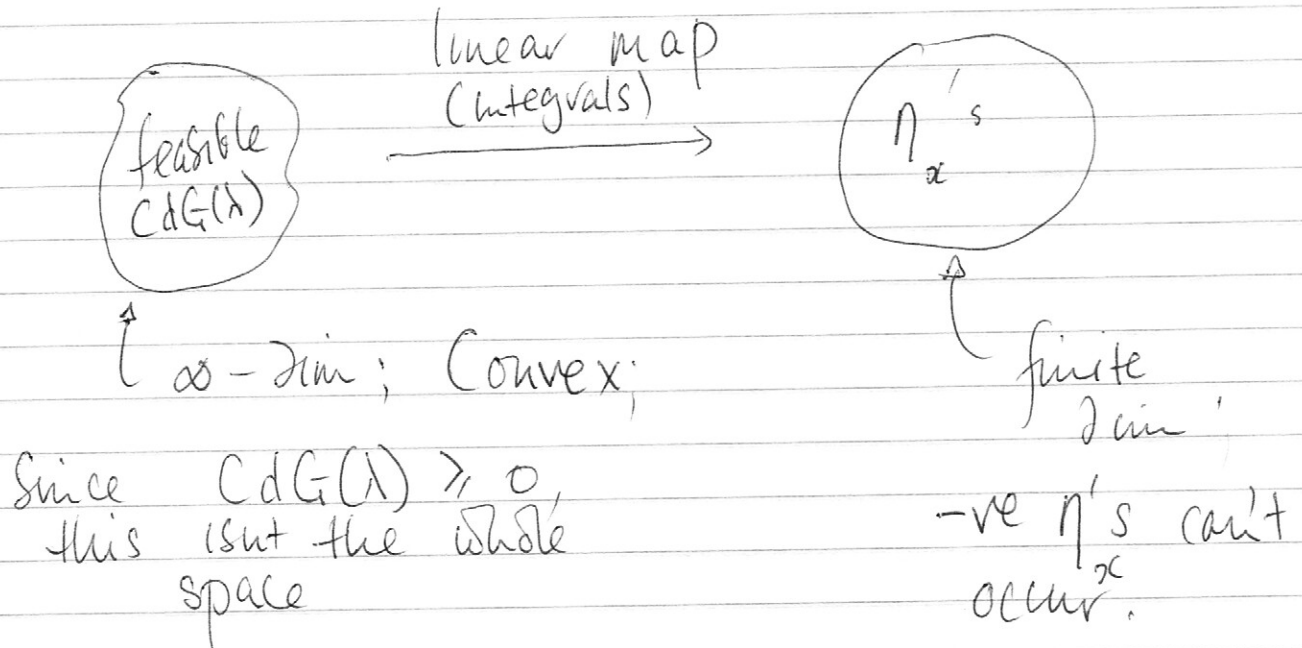
Consider $\frac{e^{-\lambda} \lambda^x}{x!}$ as a function of λ



This means we can't have

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

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

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

(17)

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

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

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

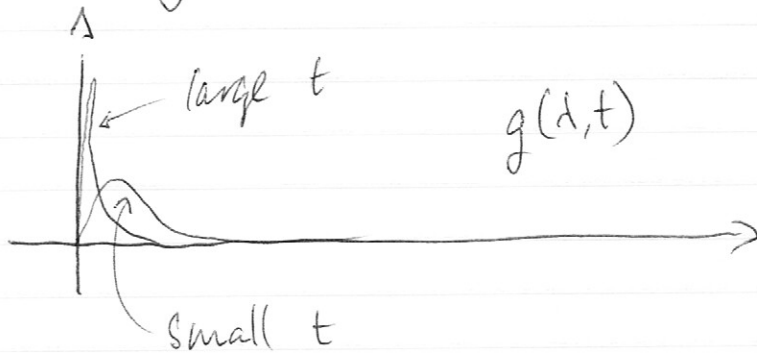
Q: how to choose μ ?

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).