BTP - Contextual Linear bandit with linear constraints

→ $\langle n, v \rangle = n^T Y$     $\|n\| = \sqrt{n^T x}$         $n =$ vector

$(d \times 1)$

$n^T$ → lamda - prob. of selecting each encoder

$(n = \pi)$                    (Try other version too!)

**Work!** $E(d \times C)$, $R(d \times 1)$

start with uniform $x$

for 1 to T :

Sample an encoder using $u_t$

get o or r by using that encoder

assign prob. It to channels with o/1

$r_t = \langle n_t, R \rangle$ → True

$C_t = \langle n_t, E Y_t \rangle$ → estimate

\# $\pi_t \cong n_t$ , $Y = tol$

Do Algorithm

**Algorithm!** $n_t (d \times 1)$, $r_t, C_t$ (Scalars)

**preprcg:**

\# $\theta^* - R$ → Normalized s.t $\Sigma R_i = 1$ (To get s)

\# $C_0 = n_0 / \|n_0\|$ , $n_0$ is safe action

$[n_{01}, n_{02} \cdots n_{0d}]$         $\Sigma n_{0i} = 1$

each row of $n_0^T E < tol$

(what if no safe action?)      find $n_0$

can we say problem          linear prog

not fisable!

# $S=1, L=1, (R=1 \text{ (doubtful !!)}) \to$ from assumptions
# Start with safe policy

$$\# \; n^T A + \hat{n}^T B + c \, \hat{n}^T \hat{n} < T$$

$$\hat{n} = x - \underset{\text{scalar}}{(n^T e)} e \quad [e \text{ and } n \text{ are }$$

$$\text{Same dimension}]$$

$$x^T R - \text{Maximize}$$

$$\sum n_i = 1 \qquad \text{Yet to figure out !}$$

---

# $J_{op} = I_{dxd} - \dfrac{1}{\|n_0\|^2} n_0 n_0^T$ $\qquad \Big| \quad \alpha_c \geq 1$

# $Sig_{opt} = \Lambda J_{op} , \quad U_{opt} = 0$

start:

# $C_{opt} = C_t - \dfrac{\langle n_t, e_0 \rangle}{\|n_0\|} C_0$

# $n_{opt} = x_t - \langle n_t, e_0 \rangle e_0$

# $Sig_{opt} \mathrel{+}= n_{opt} \times (n_{opt})^T$

# $U_{opt} \mathrel{+}= C_{opt} \times n_{opt}$

# $M_{opt} = (Sig_{opt}^{-1}) \times U_{opt}$

# $\tilde{C}_{\pi t} < T \quad$ (below are simplification steps)

$$\Rightarrow \frac{C_0 C_0}{\|n_0\|} z^t + \Big( z - e_0 z^T e_0 \Big)^T M_{opt} + \underset{\| x_t \|}{\alpha_c \beta_t}$$

$$\Downarrow$$

$$\big( z^T - e_0^T z e_0^T \big) M_{opt}$$

$$\implies \frac{C_0}{\|n_0\|} z^T e_0 + z^T M_{opt} + e_0^T z (e_0^T M_{opt}) + \underset{|x|}{\alpha_c \beta_t} \|n_t\|$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$

$$e_0^T z \qquad M_{opt}^T z \qquad (e_0^T M_{opt}) e_0^T z$$

$$\Rightarrow \left( \frac{c_0}{\|n_0\|} c_0^T + M_{opt}^T + c_0^T n_{opt} c_0^T \right) z + \alpha_c \beta_t \|x_{t-1}\|$$

$$\leq T \to \text{bound}$$

approx

$R^T z \to$ maximize

$\leq z_i \geq 1 \to eq$

## Notes on Final Implementation:

1) paper assumes

$$r_t = \langle n_t, \theta_* \rangle + \xi_t^r \quad \text{but in our case wk}$$

$\theta_* = R_{ate}$ & $\xi_t^r = 0$ precisely so avoided

the calculations of $\theta_t$

2) To get a safe action the min_tol should be higher than the min_tol for previous algorithms

3) In cases like $n_0$ or $n_t$ close to $[1,0,0]$ the $M u_{opt}$ is not possible (singular) or blowing up. hence performing pseudo inverse

4) Final eqn. $Az + C \leq T$

$$Az \leq T - C$$

here $T = 0.2 | 0.3$,

$C = \alpha_c \beta_t \|n_{t-1}\| \simeq \alpha_c \cdot 3.2$

and coeff A are mostly +ve or slightly -ve

So for $\alpha_c \geq 1$ as given in paper we are getting an infeasible equation

+ For our problem best $dc = 0.001$ but any value below $0.01$ is working

5) The rate is decreasing starting from a high value to adjust error, In prev. algorithms rate used to increase from a lower value may be due to different starting points (no, uniform)