

Last time: 1) More Fisher info.

2) Cramér-Rao bound

3) Differential correlations

fW: 1) How to estimate θ from x, y

2) Fisher information lower bound

Information Theory

Def: The surprise of an outcome k w/ probability $P(k)$

is:

$$u_k = -\log_2 P(k) \text{ (in bits)}$$

Example 1: Weighted coin flip, prob. of heads $P(H)$

$$P(H) = \frac{1}{2}, \quad u_H = -\log_2 \frac{1}{2} = 1 \text{ bit}$$

$$P(T) = 1 \quad u_T = 0 \text{ bits}$$

$$P(H) = \frac{1}{4} \quad u_H = 2 \text{ bits}$$

Example 2: Two coins, $P(H) = \frac{1}{2}$, $P(T|H) = \frac{1}{4}$

$$u_{HH} = 2 \text{ bits}$$

Surprise is additive for independent events.

Def: The entropy of a r.v. X that takes values x_k w/ probabilities $P(x_k)$ is the avg. surprise:

$$H(X) = - \sum_k P(x_k) \log_2 P(x_k)$$

Example:

$$\text{If } P(H) = P(T) = \frac{1}{2}$$

$$H = -P(H) \log_2 P(H) - P(T) \log_2 P(T) = 1 \text{ bit}$$

$$P(H) = 1, P(T) = 0 \Rightarrow H = 0 \text{ bits}$$

$$P(H) = \frac{3}{4}, P(T) = \frac{1}{4} \Rightarrow H = 0.81 \text{ bits}$$

For discrete prob. dists., uniform dist. has highest entropy.

If n outcomes, then $P(x_k) = \frac{1}{n}$, $k = 1 \dots n$

$$H(X) = - \sum_{k=1}^n \frac{1}{n} \log_2 \frac{1}{n} = - \log_2 \frac{1}{n} = \log_2 n$$

Def: The mutual information between s and r is:

$$\begin{aligned} I(r,s) &= H(r) - H(r|s) = H(s) - H(s|r) \\ &= H(r) + H(s) - H(r,s) \end{aligned}$$

where $H(r|s)$ is the conditional entropy

$$\sum_j P(s_j) H(r|s=s_j) = - \sum_j P(s_j) \sum_k P(r_k|s_j) \log_2 P(r_k|s_j)$$

(assuming that r takes values $\{r_k\}$ and s takes $\{s_j\}$)

Intuition using Kullback-Leibler divergence:

Given two prob dists. $P(x)$, $Q(x)$,

$$D_{KL}(P \parallel Q) = \sum_k P(x_k) \log_2 \frac{P(x_k)}{Q(x_k)}$$

(need $Q(x_k) = 0 \Rightarrow P(x_k) = 0$)

Properties: 1) $D_{KL}(P \parallel Q) \underset{>0}{\searrow} 0$ if $P(x_k) = Q(x_k) \quad \forall k$

2) $D_{KL} \geq 0$

3) $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$

$$\text{Claim: } I(r,s) = D_{KL}(P(r,s) \parallel P(r)P(s))$$

$$= \sum_{j,k} P(r_k, s_j) \log_2 \frac{P(r_k, s_j)}{P(r_k)P(s_j)}$$

$$\text{Note: } P(r_k, s_j) = P(r_k | s_j) P(s_j)$$

$$= \sum_j P(s_j) \sum_k P(r_k | s_j) \left| \log_2 \frac{P(r_k | s_j)}{P(r_k)} \right|$$

$$= \sum_j P(s_j) \sum_k P(r_k | s_j) \log_2 P(r_k | s_j) \quad (= -H(r|s))$$

$$- \sum_j P(s_j) \sum_k P(r_k | s_j) \log_2 P(r_k) \quad (= H(r))$$

$$\sum_k P(r_k)$$

$$= H(r) - H(r|s) = I(r,s)$$

Same argument for $H(s) - H(s|r)$

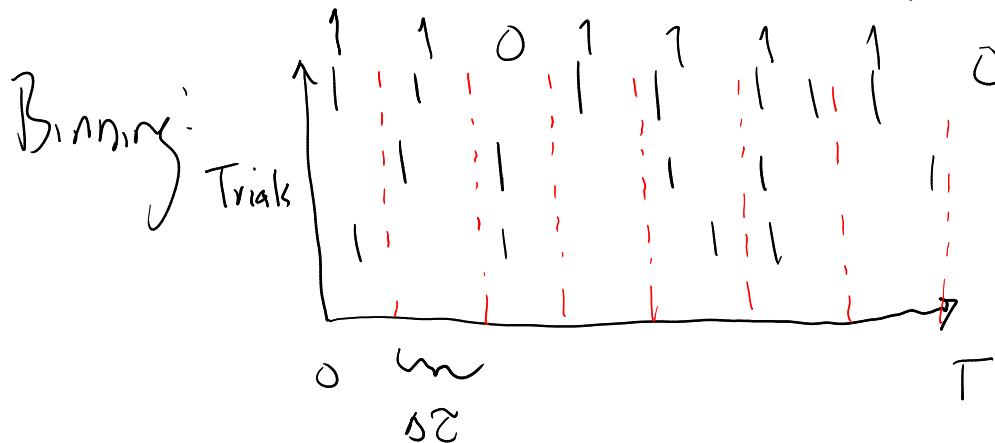
Conclusion: Mutual info. is KL-divergence b/w joint & independent dists. of r, s.

Application 1: Entropy of a spike train

Given spike times $P(\{t_k\} | s)$, how much info. about s is present?

Problem: Infinite-dimensional space of $\{t_k\}$.

Possibilities: Consider statistics (e.g. spike count)



$$N = T/\Delta \tau \text{ bins}$$

If 0/1, 2^N possible binary "words"

Hard to estimate for large T .

Depends on $\Delta \tau$.

"Direct method" (Strong, Kobele, de Ruyter van Steveninck, Bialek, PRL 1997)

Example (Strong et al. 1997):

If $P(x_k)$ is uniform over n outcomes, $H = \log_2 n$

and 2 random observations are the same w/ prob. $1/n$.

Can use the prob. of "coincidence" P_c to

estimate $H \approx -\log_2 P_c$

$$\text{In general, } P_c = \sum_k P(x_k)^2 = E[P(x_k)]$$

$$H = - \sum_k P(x_k) \log_2 P(x_k) = -E[\log_2 P(x_k)]$$

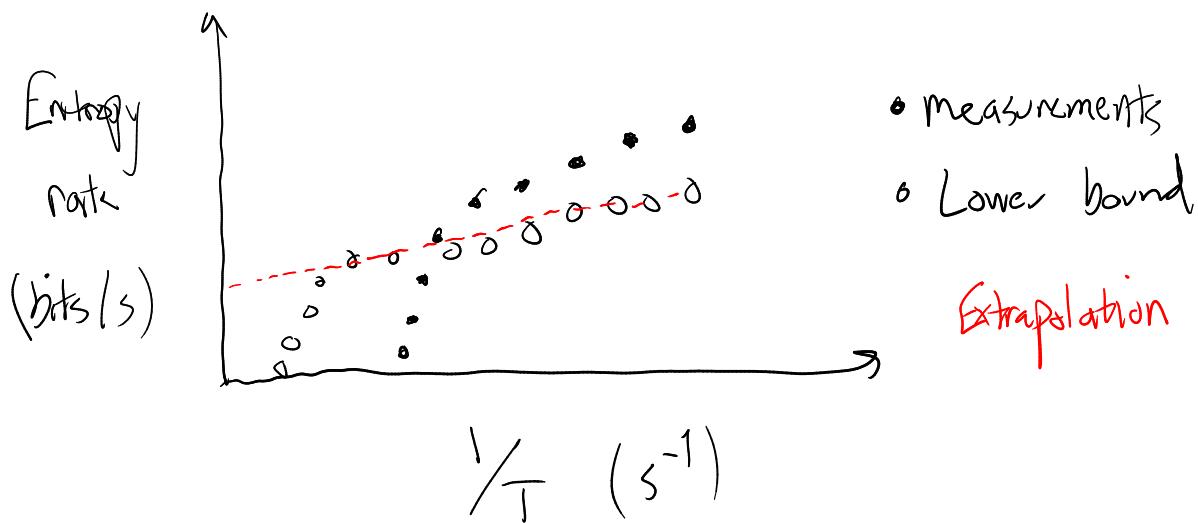
$$\geq -\log_2 E[P(x_k)] = -\log_2 P_c$$

For spike trains, assume words w/ same # of spikes S equally probable.

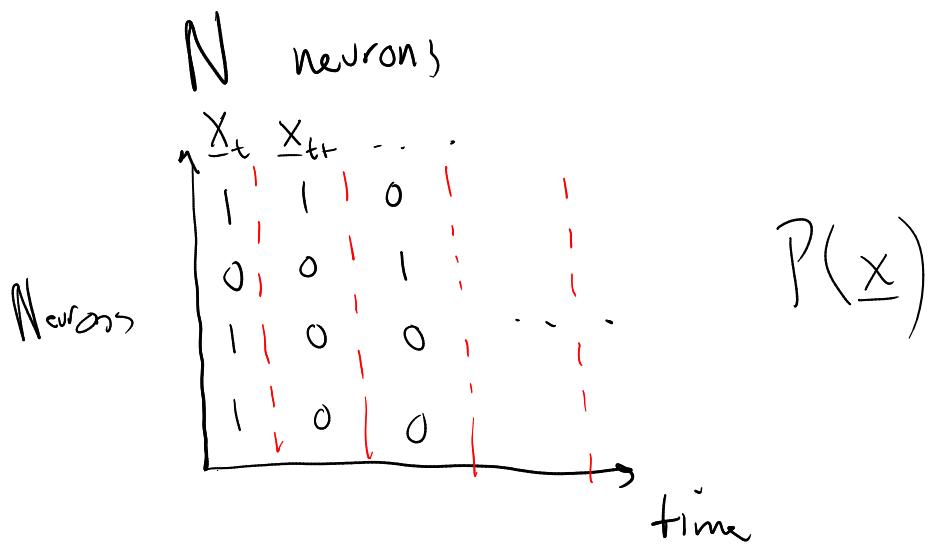
$$H_{LB} = - \sum_S P(S) \log_2 \left[P(S) \frac{2N_c(S)}{N_{obs}(S)[N_{obs}(S)-1]} \right]$$

$N_c(S)$: # of coincidences w/ S spikes

$N_{obs}(S)$: # total occurrences of words w/ S spikes



Application 2: Maximum entropy models



How much structure in this population code can be explained by constraints?

Idea: Find maximum entropy distribution consistent w/
 2^N a set of constraints. Which constraints?

i) $\sum_{k=1}^{2^N} P(\underline{x}_k) = 1$ ii) $\sum_k P(\underline{x}_k) x_i = \mu_i = E[x_i]$
 (avg. rate of neuron i).

$$\sum_k P(x_k) x_i x_j = E[x_i x_j] = c_{ij} \quad (\text{correlations})$$

$$\text{let } P(x_k) = P_k$$

Maximize $-\sum_k P_k \log_2 P_k$ subject to these constraints

Method of Lagrange multipliers:

$$\begin{aligned} \frac{\partial}{\partial P_k} & \left\{ -\sum_k P_k \log_2 P_k + \lambda \left(\sum_k P_k - 1 \right) \right. \\ & + \sum_{i=1}^N h_i \left(\sum_k P_k x_i - u_i \right) \\ & \left. + \sum_{i,j=1}^N J_{ij} \left(\sum_k P_k x_i x_j - c_{ij} \right) \right\} = 0 \end{aligned}$$

$$= -\log P_k - 1 + \lambda + \sum_i h_i x_i + \sum_{ij} J_{ij} x_i x_j$$

$$P(x_k) = \exp \left(\lambda - 1 + \sum_i h_i x_i + \sum_{ij} J_{ij} x_i x_j \right)$$

$$= \frac{1}{Z} \exp \left(\sum_i h_i x_i + \sum_{ij} J_{ij} x_i x_j \right)$$

normalization first-order
constraints Second-order
 "Boltzmann distribution"

If no constraints on h_i, c_{ij}

$$P(x_k) = \frac{1}{Z}, \text{ uniform dist.}$$

If only on h_i

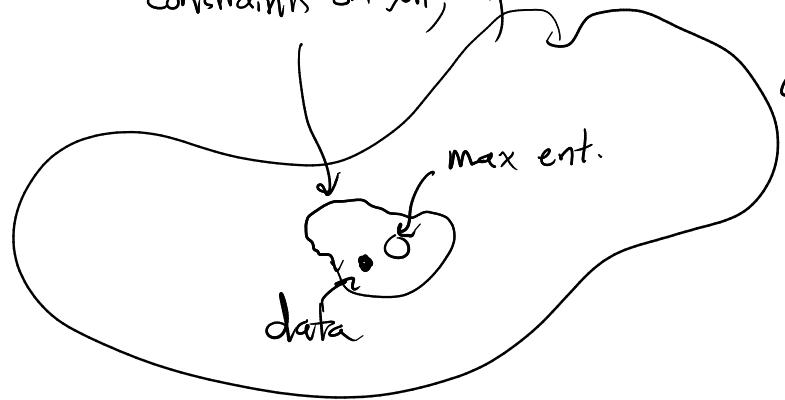
$$P(x_k) = \frac{1}{Z} \exp \left(\sum_i h_i x_i \right) = \frac{1}{Z} \prod_i e^{h_i x_i}$$

(independent)

All dists $P(x)$

constraints on h_i, c_{ij}

max ent.



dists.
Satisfying constraints
on h_i

If H_L is entropy of max-ent. dist satisfying L^{th} order
constraint,

$$H_1 > H_2 > \dots > H_N$$

Schneidman et al Nature 2006