

Introduction to Artificial Intelligence

Lecture 14 – Information gathering

CS/CNS/EE 154

Andreas Krause

Announcements

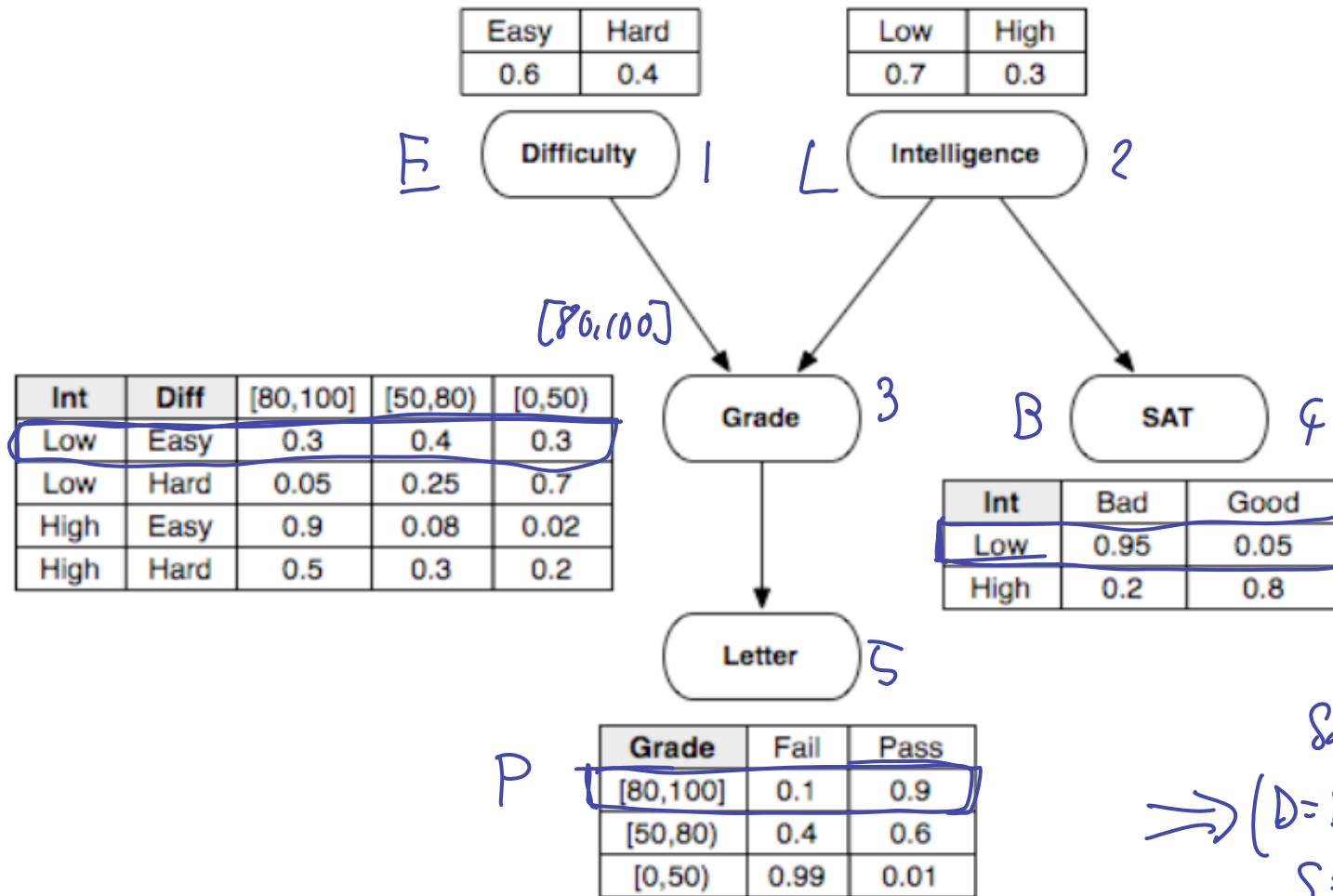
- Homework 2 due today
- Homework 3 out later this week
- Final project due December 1
- Code released on Monday (Nov 15)

- Note on midterm grades (Avian Asker)

Sampling based inference

- So far: deterministic inference techniques
 - Variable elimination
 - (Loopy) belief propagation
- Will now introduce stochastic approximations
 - Algorithms that “randomize” to compute expectations
 - In contrast to the deterministic methods, guaranteed to converge to right answer (if wait looong enough..)
 - More exact, but slower than deterministic variants

Forward sampling from a BN



Sample
 $\Rightarrow (D=E, I=L, G=[10, 100],$
 $S=B, L=P)$

Rejection sampling

- Collect samples over all variables

$$\hat{P}(\mathbf{X}_A = \mathbf{x}_A \mid \mathbf{X}_B = \mathbf{x}_B) \approx \frac{\text{Count}(\mathbf{x}_A, \mathbf{x}_B)}{\text{Count}(\mathbf{x}_B)}$$

- Throw away samples that disagree with \mathbf{x}_B
- Can be problematic if $P(\mathbf{X}_B = \mathbf{x}_B)$ is rare event

Sample complexity for probability estimates

- Absolute error:

$$\text{Prob}\left(|\hat{P}(\mathbf{x}) - P(\mathbf{x})| > \varepsilon\right) \leq 2 \exp(-2N\varepsilon^2)$$

- Relative error:

$$\text{Prob}\left(\hat{P}(\mathbf{x}) < (1 + \varepsilon)P(\mathbf{x})\right) \leq 2 \exp(-NP(\mathbf{x})\varepsilon^2/3)$$

if $P(\mathbf{x})$ exponentially small, need N exponentially large

Sampling from rare events

- Estimating conditional probabilities $P(X_A \mid \mathbf{X}_B = \mathbf{x}_B)$ using rejection sampling is hard!
 - The more observations, the unlikelier $P(\mathbf{X}_B = \mathbf{x}_B)$ becomes
- Want to directly sample from posterior distribution!

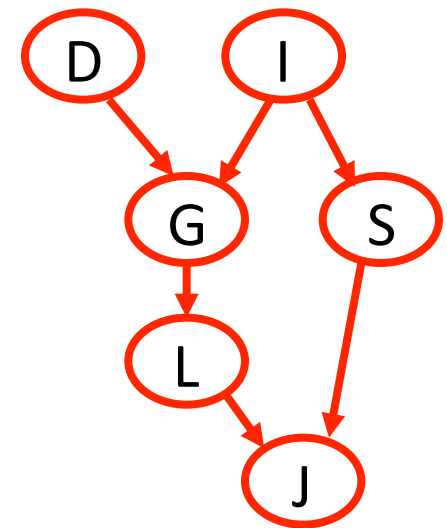
Gibbs sampling

- Start with initial assignment $\mathbf{x}^{(0)}$ to all variables
- For $t = 1$ to ∞ do
 - Set $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)}$
 - For each variable X_i
 - Set $\mathbf{v}_i =$ values of all $\mathbf{x}^{(t)}$ except x_i
 - Sample $x_i^{(t)}$ from $P(X_i \mid \mathbf{v}_i)$
- For large enough t , sampling distribution will be “close” to true posterior distribution!
- **Key challenge:** Computing conditional distributions $P(X_i \mid \mathbf{v}_i)$

Gibbs Sampling

Gibbs sampling $P(D, I, G, S, L \mid J = 1)$

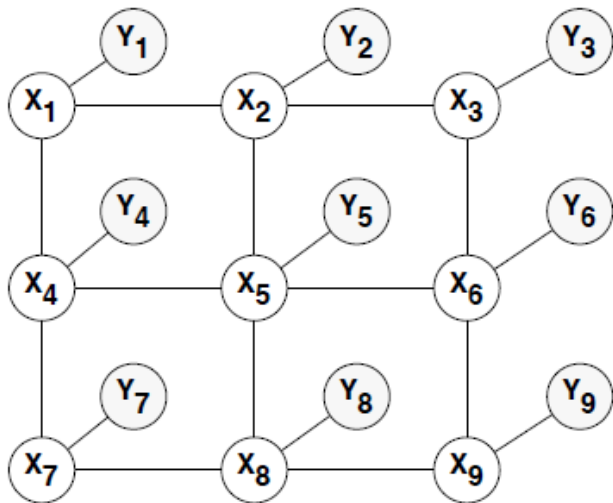
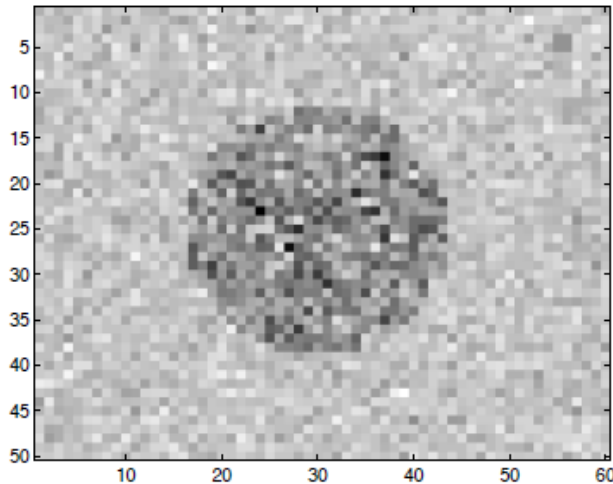
Iter	D	I	G	S	L	J
1	0	1	1	0	0	1
2	* 0	1	1	0	1	1
3	0	0	1	0	0	1
4	0	0	1	0	1	1
...						



$P(G=1) \approx \frac{3}{4}$

$$\begin{aligned}
 * & \sim P(D \mid I=1, G=1, S=0, L=0, J=1) \\
 &= \frac{1}{2} P(D, I=1, G=1, S=0, L=0, J=1) = \\
 &= \frac{P(D) P(I=1) P(G=1 \mid D, I=1) P(S=0 \mid I=1) P(L=0 \mid G=1) P(J=1 \mid L=0, S=0)}{\sum_d P(D=d) P(I=1) P(G=1 \mid D=d, I=1) P(S=0 \mid I=1) P(L=0 \mid G=1) P(J=1 \mid L=0, S=0)} \\
 &= \frac{1}{2} P(D) P(G=1 \mid D, I=1)
 \end{aligned}$$

Example: (Simple) image segmentation



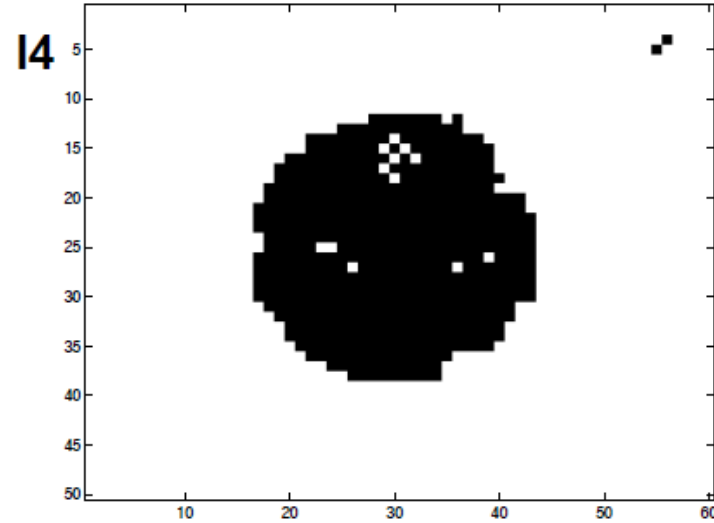
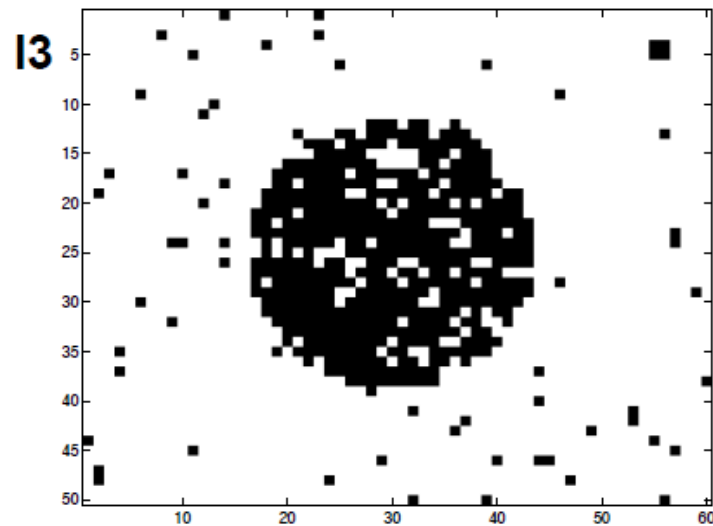
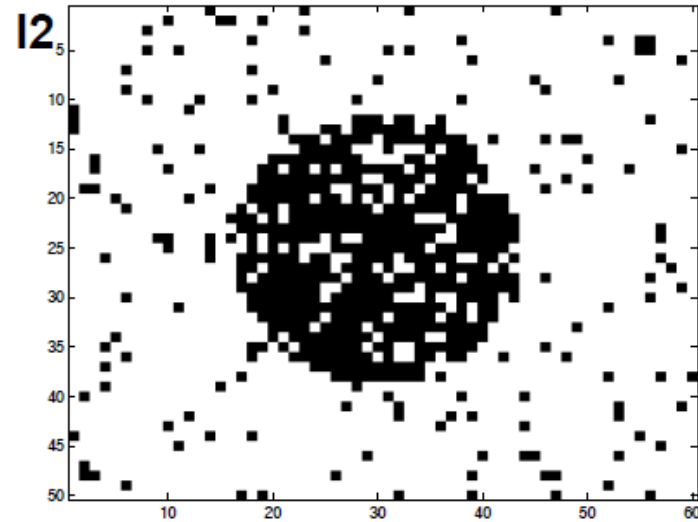
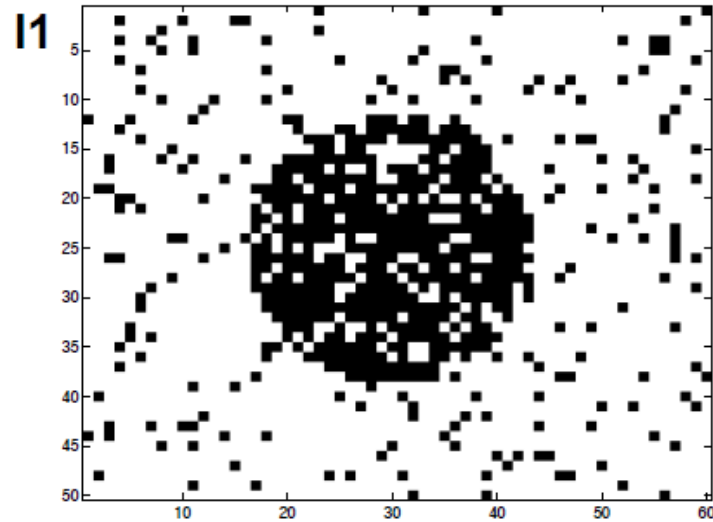
$$P(x) = \frac{1}{Z} \prod_i \Phi(x_i) \prod_{(j,k) \in E} \Psi(x_j, x_k)$$

$$\Phi(x_i) = \exp \left\{ -\frac{(y_i - \mu_{x_i})^2}{2\sigma_{x_i}^2} \right\}$$

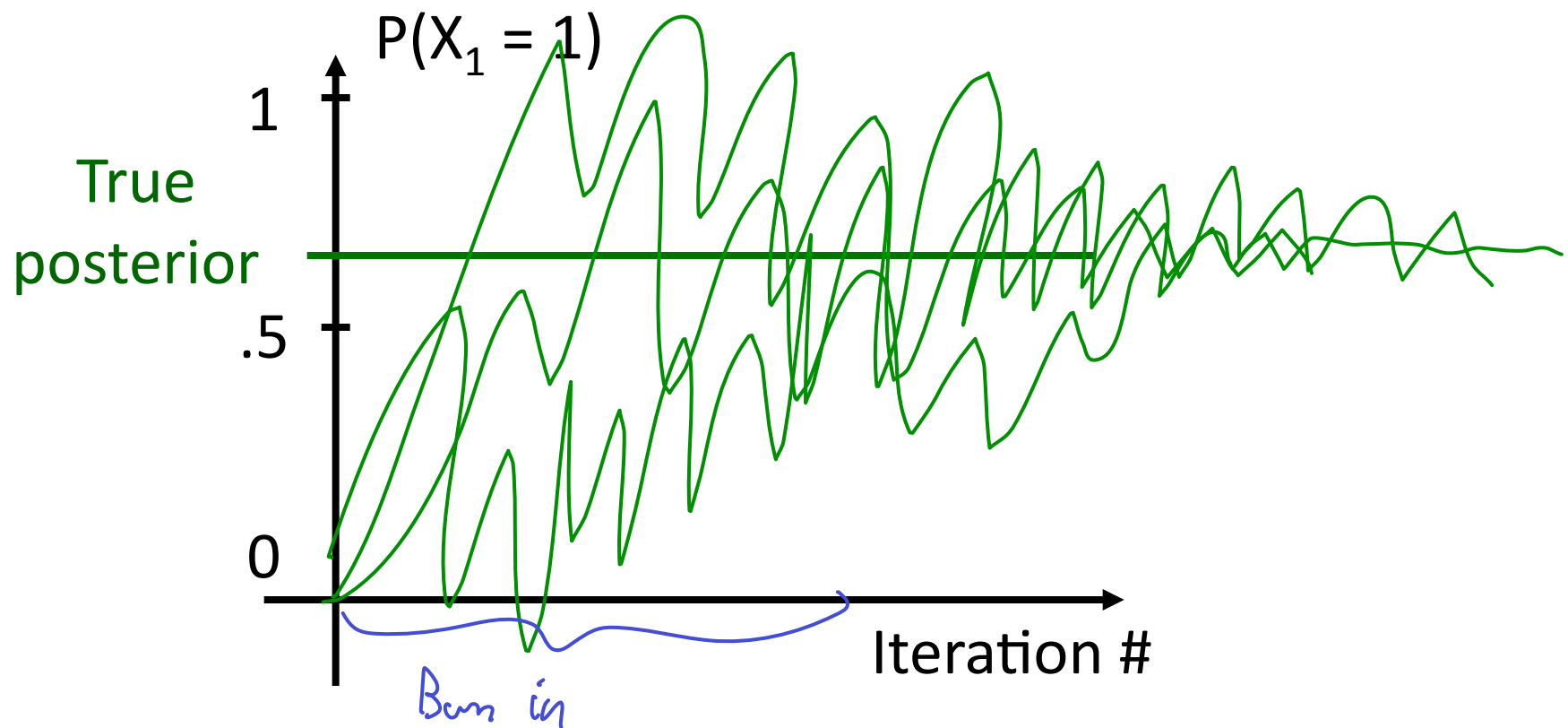
$$\Psi(x_i, x_j) = \exp \{ -\beta(x_i - x_j)^2 \}$$

[see Singh '08]

Gibbs Sampling iterations



Convergence of Gibbs Sampling



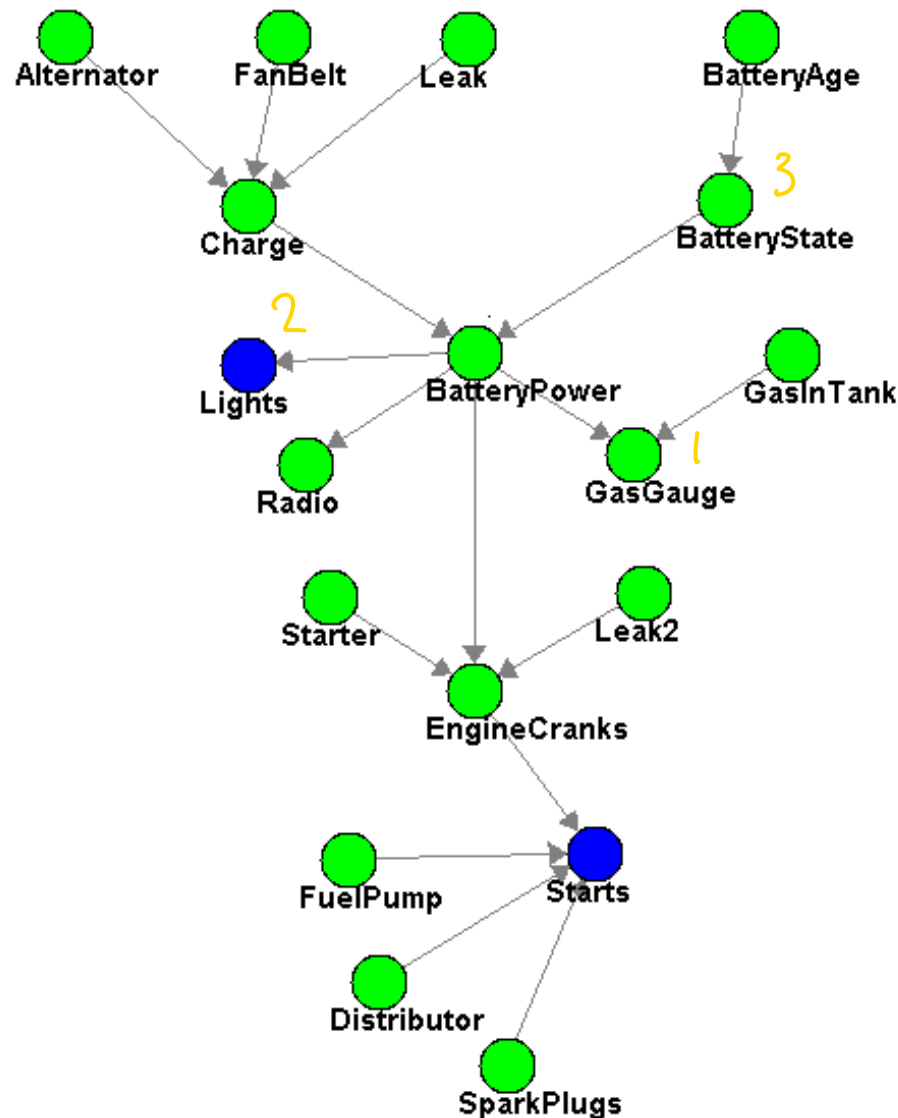
Summary: Inference

- For tree-structured Bayes nets, can compute exact marginals
 - Variable elimination
 - Belief propagation (efficiently computes all marginals)
- For loopy networks, can use approximate inference
 - Loopy belief propagation (may not converge)
 - Gibbs sampling (will converge, but may take long time)

Information gathering

- So far:
 - Bayesian networks for quantifying uncertainty in real world environments
 - Exact and approximate algorithms for inference in Bayesian networks (e.g., compute $P(\text{Pit} \mid \text{Breezes})$)
- Now:
 - Selecting most “informative” variables for making effective predictions / decisions

Why does my car not start?



- Selectively run tests to diagnose cause of failure

Clinical diagnosis?

- Patient either healthy or ill
- Can choose to treat or not treat

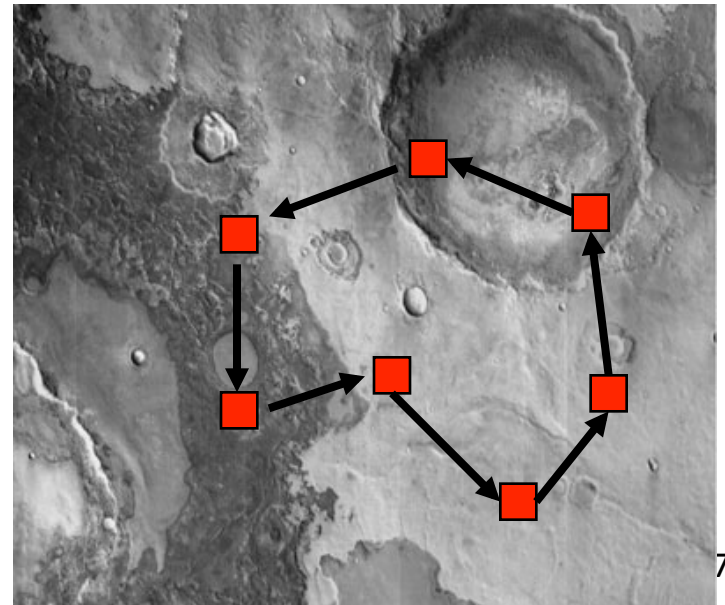
	<i>healthy</i>	<i>ill</i>
Treatment	-\$-\$	\$
No treatment	0	-\$\$\$

- Only know distribution $P(\text{ill} \mid \text{observations})$
- Can perform costly medical tests to reveal aspects of the condition
- **Which tests should we perform to most cost-effectively diagnose?**

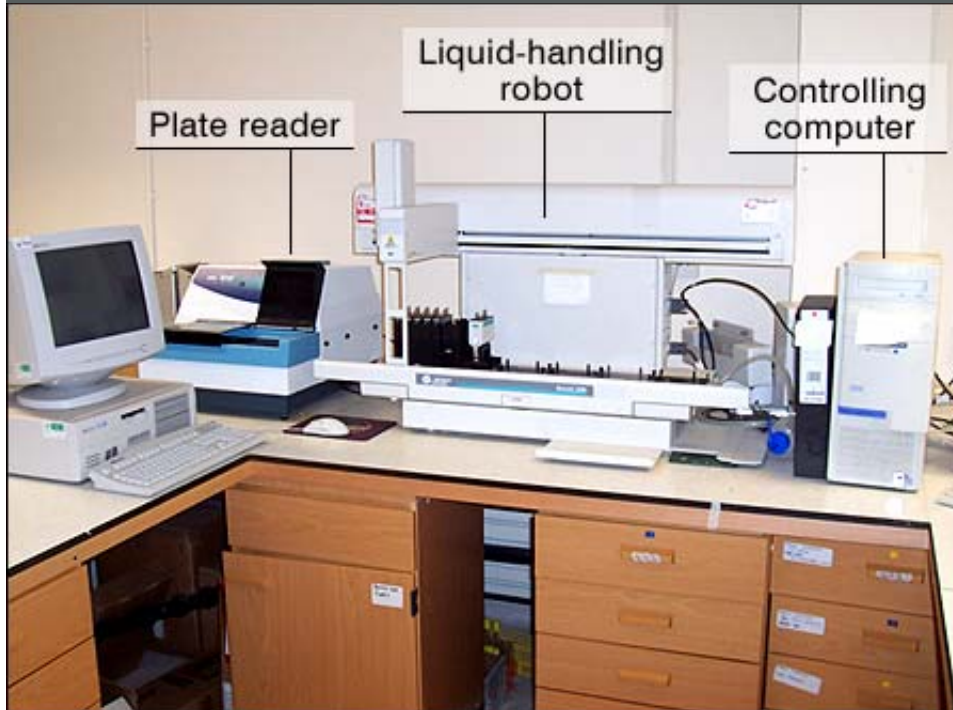
Autonomous robotic exploration



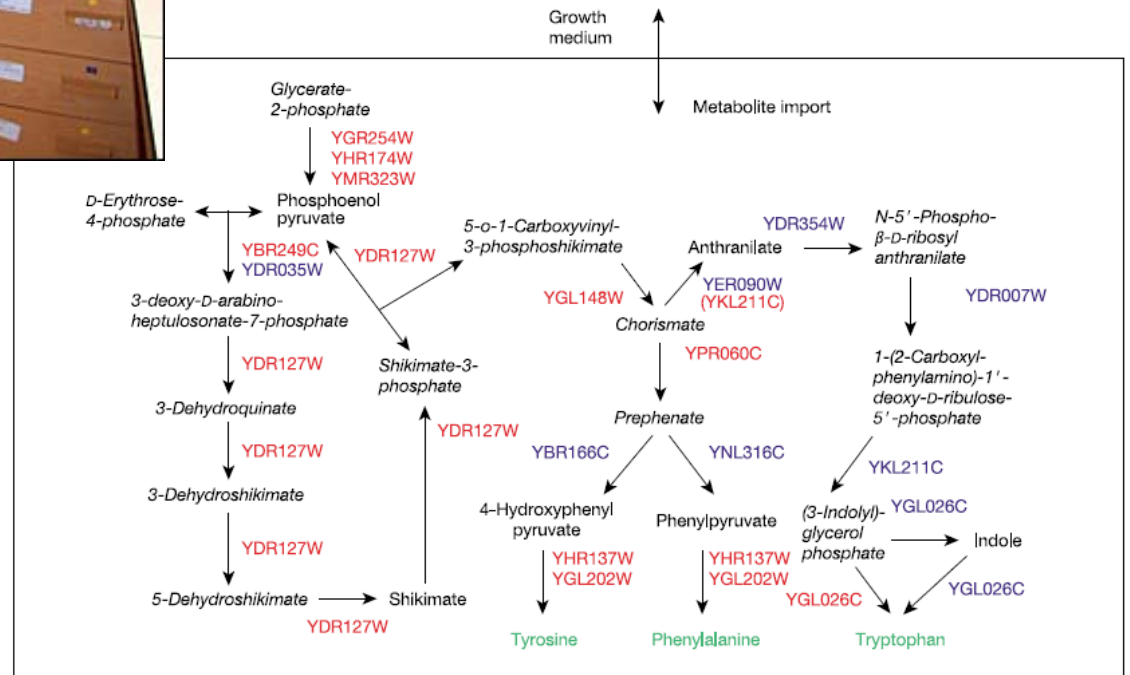
- Limited time for measurements
- Limited capacity for rock samples
- **Need optimized information gathering!**



A robot scientist

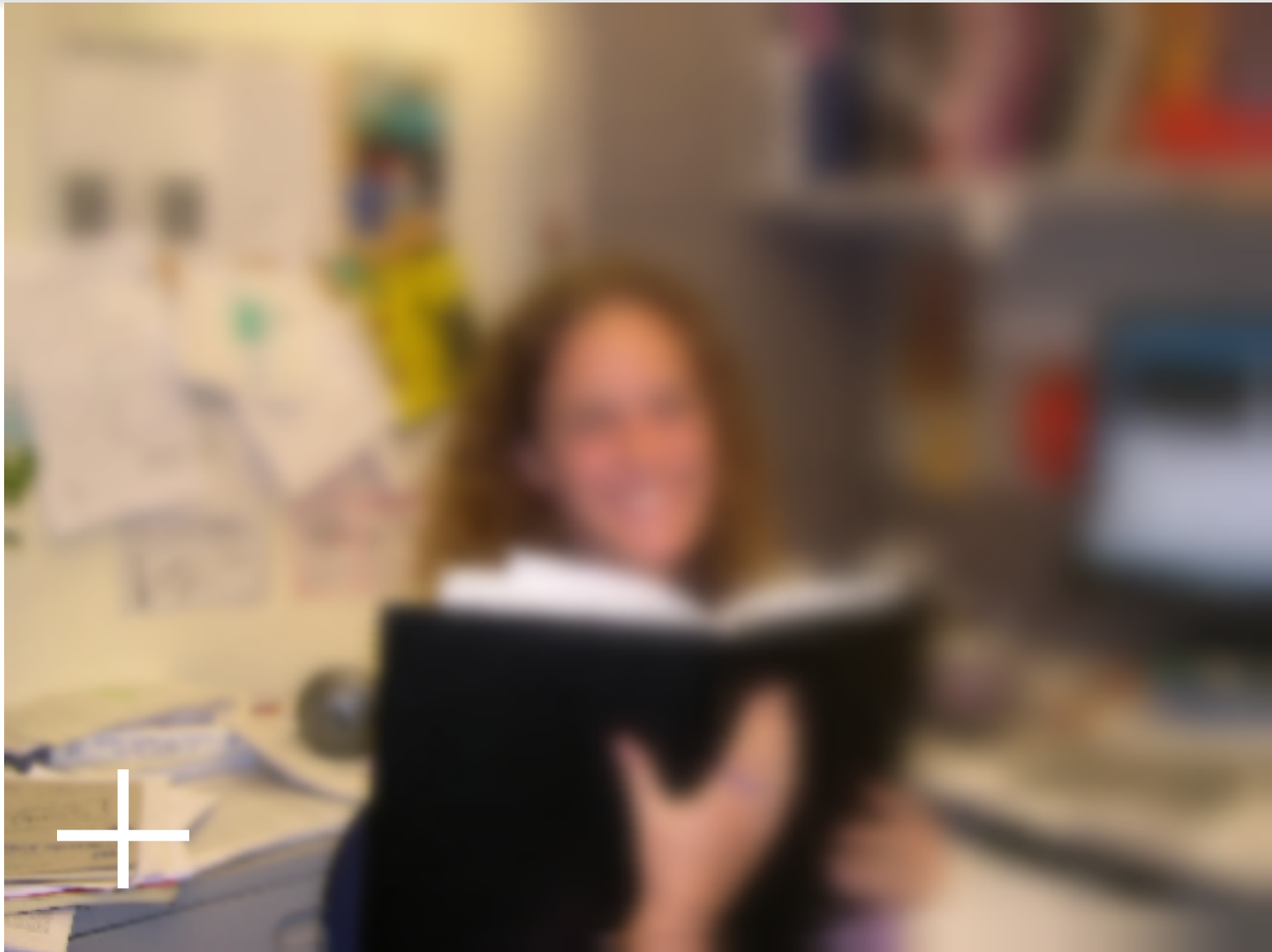


King et al, Nature '04



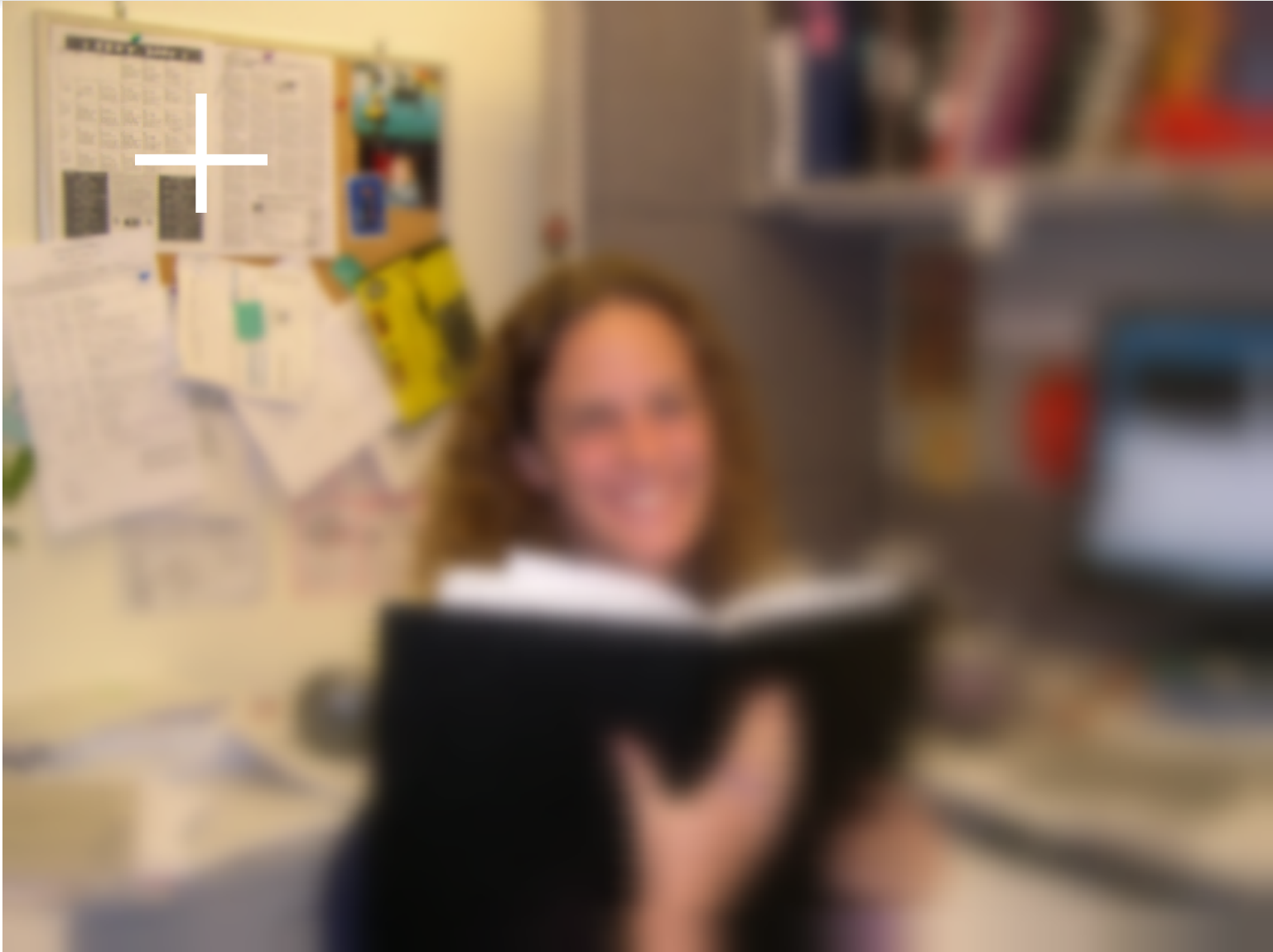
How do people gather information?

[Renninger et al, NIPS '04]



How do people gather information?

[Renninger et al, NIPS '04]



How do people gather information?

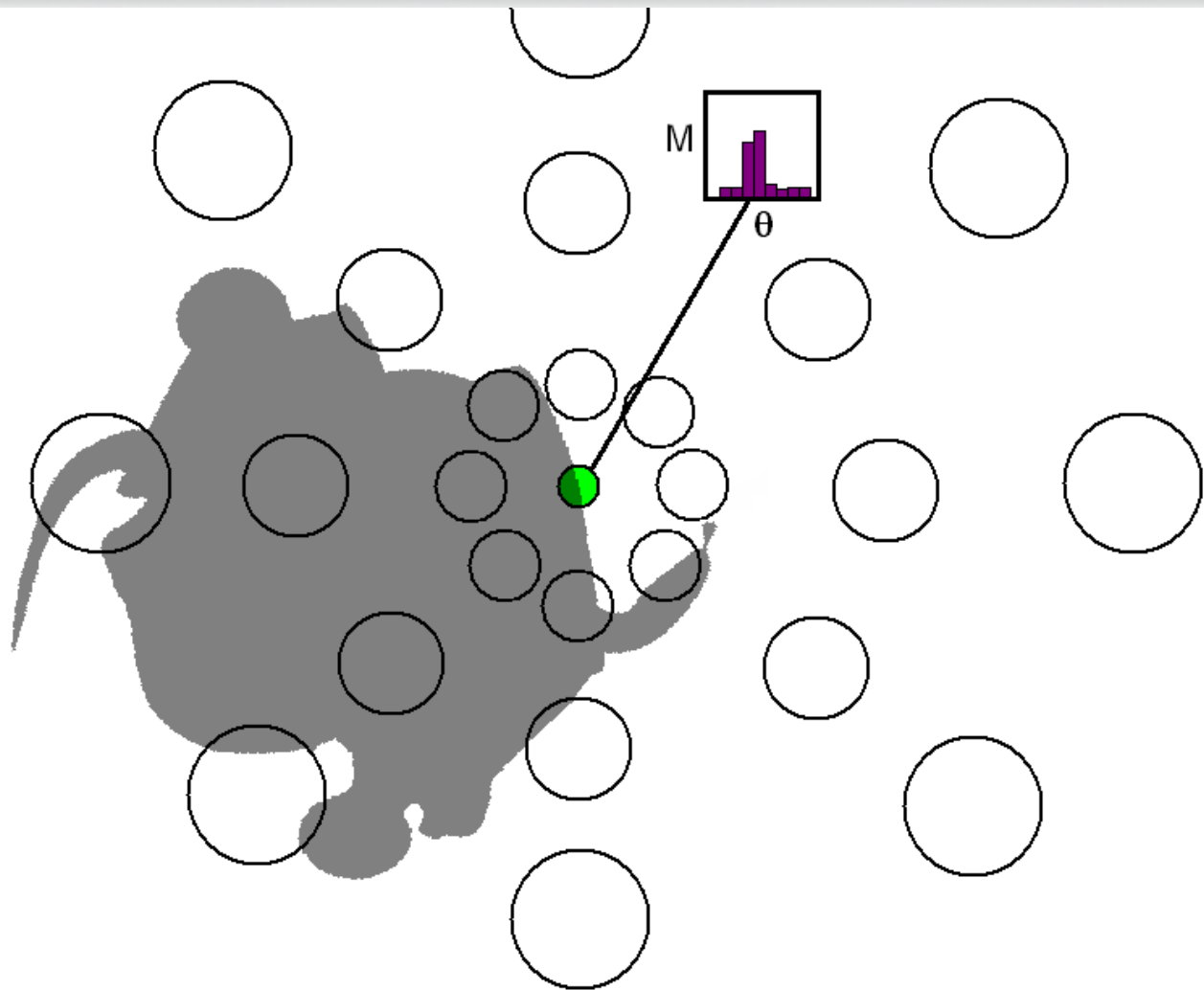
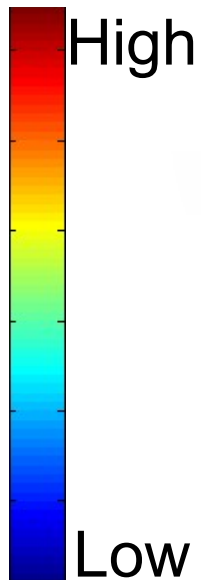
[Renninger et al, NIPS '04]



How do people gather information?

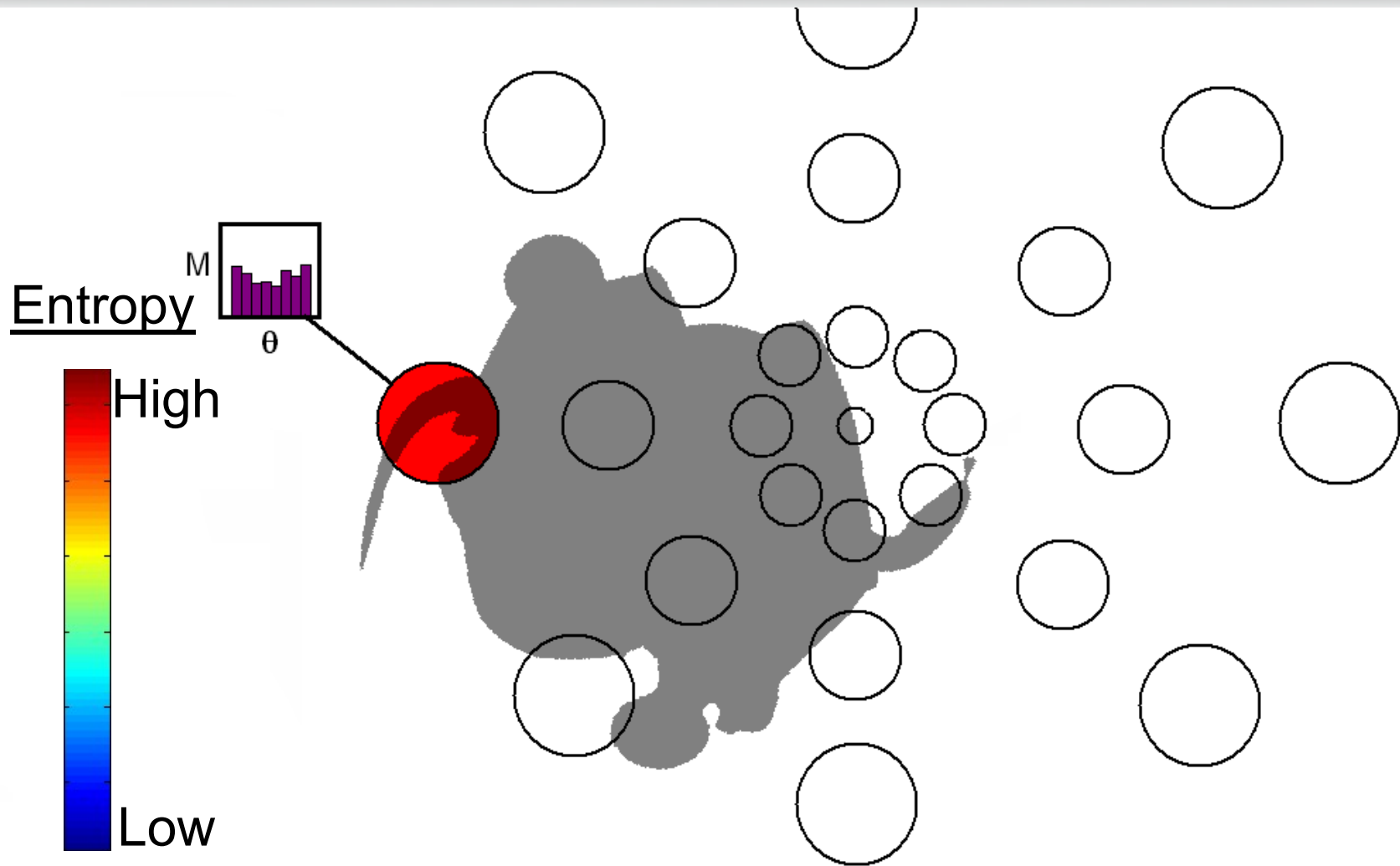
[Renninger et al, NIPS '04]

Entropy



How do people gather information?

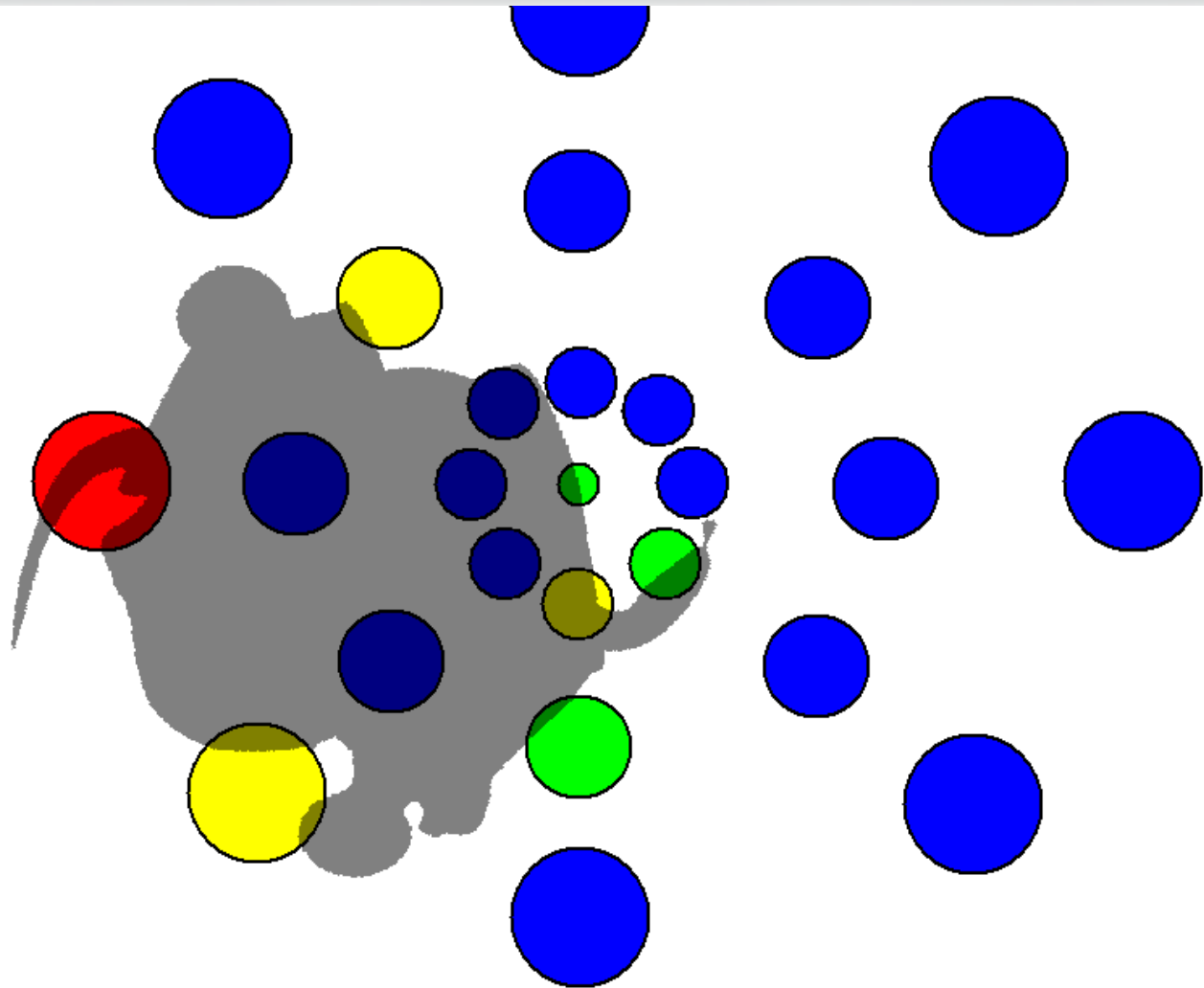
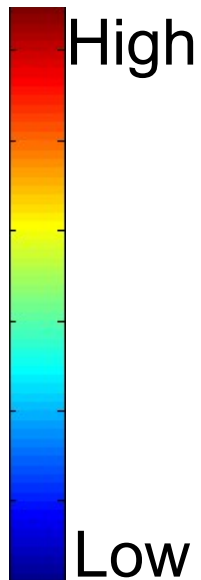
[Renninger et al, NIPS '04]



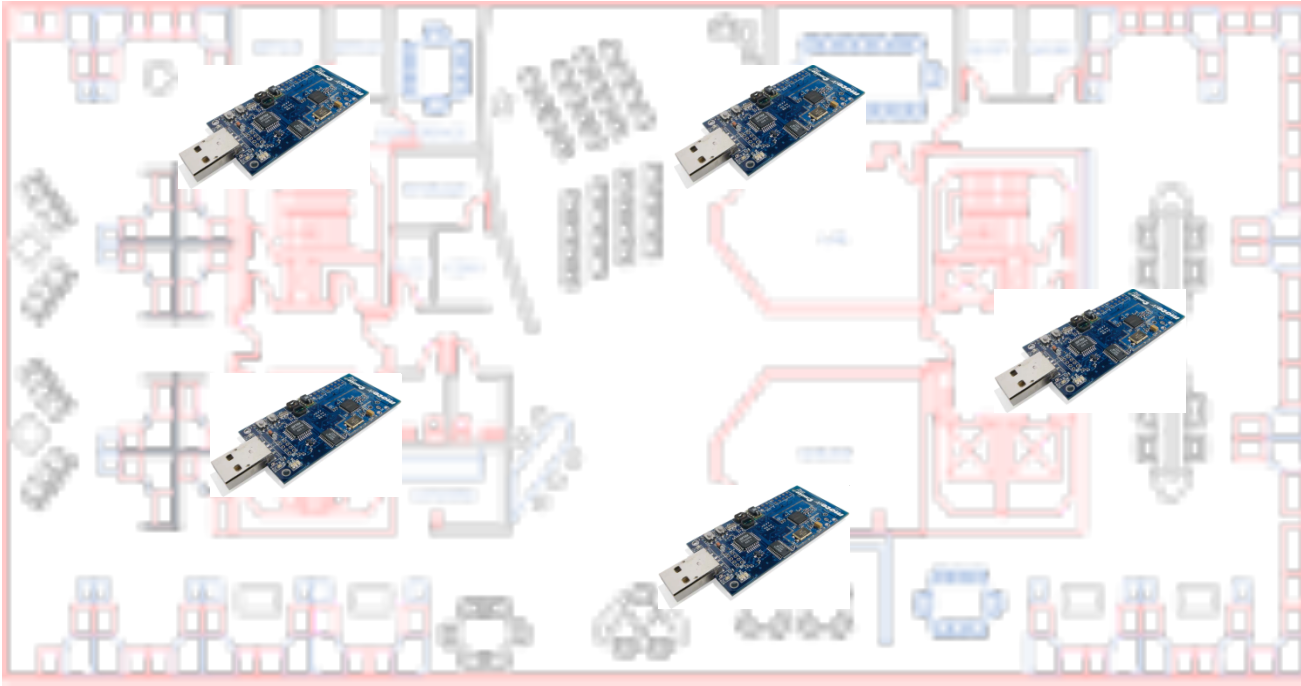
How do people gather information?

[Renninger et al, NIPS '04]

Entropy

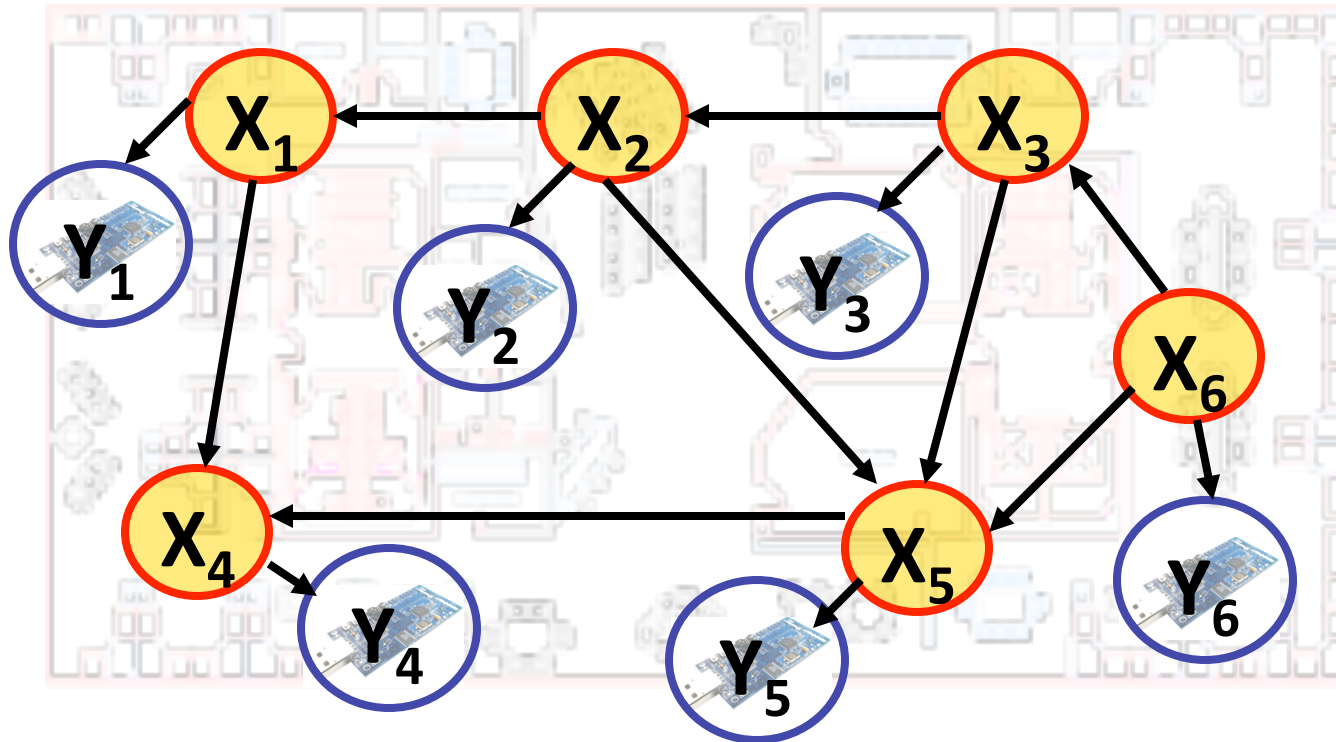


Running example: Detecting fires



Want to place sensors to detect fires in buildings

Monitoring using Bayesian Networks



X_s : temperature
at location s

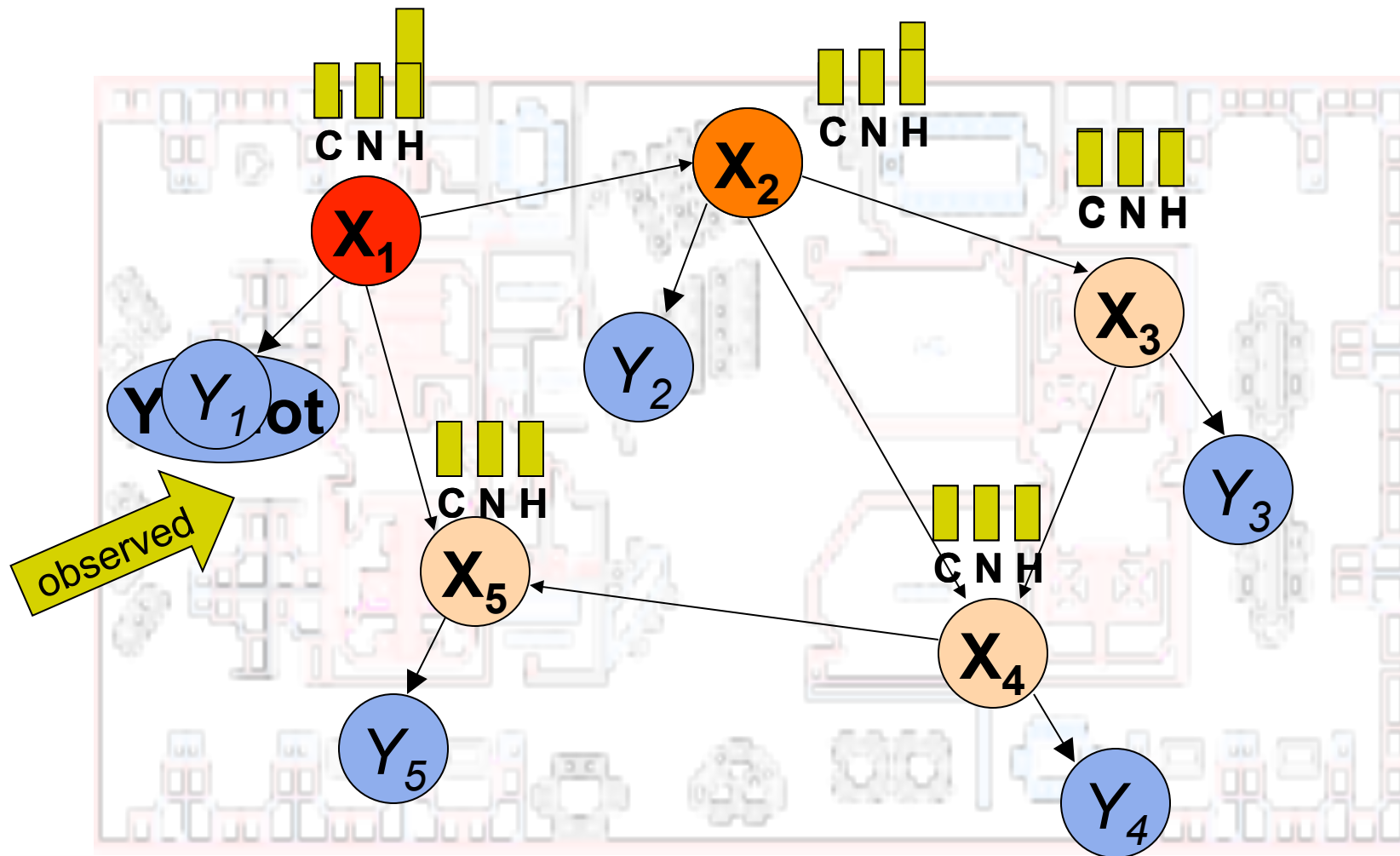
Y_s : sensor value
at location s

$$Y_s = X_s + \text{noise}$$

Joint probability distribution

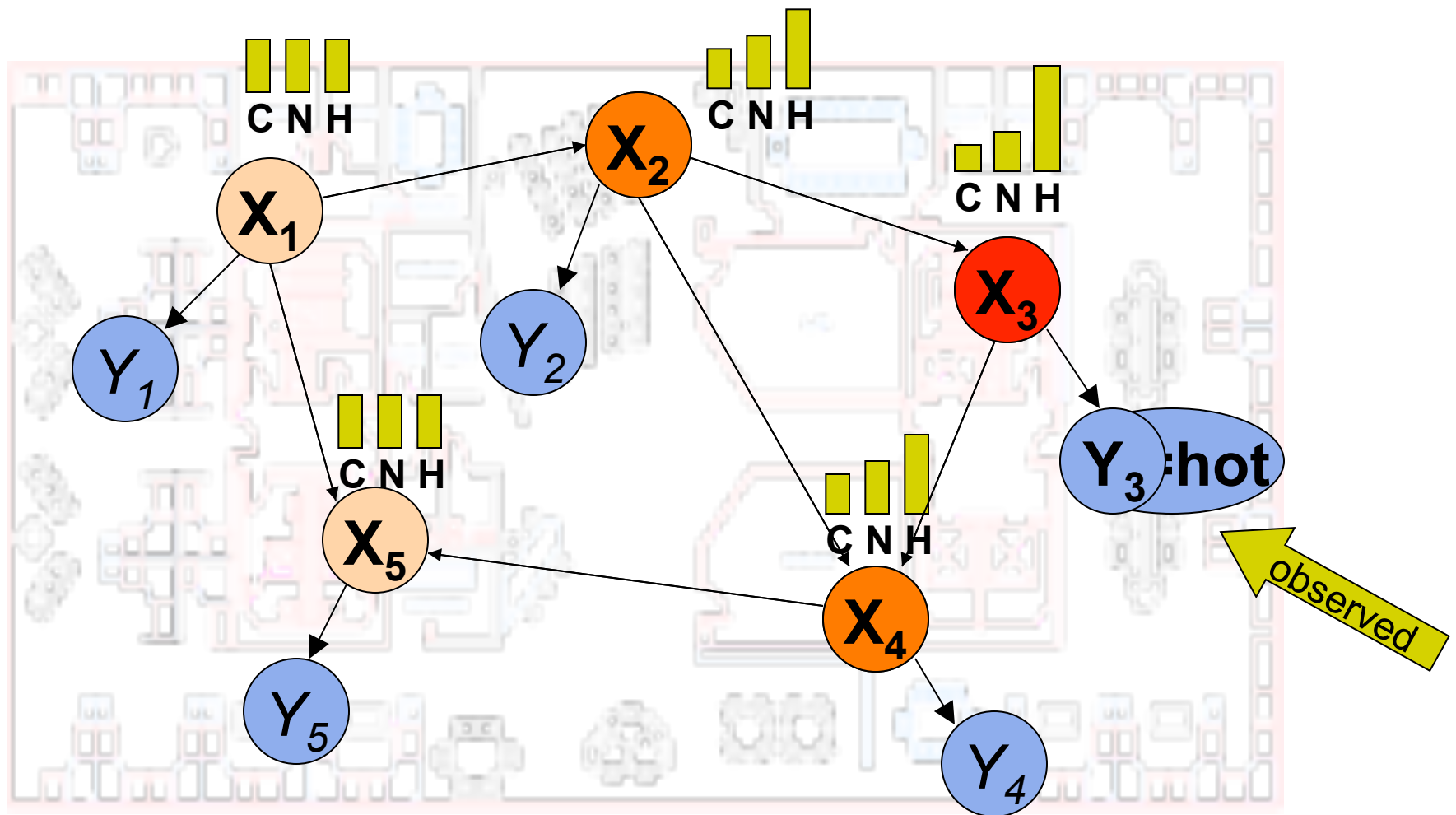
$$P(X_1, \dots, X_n, Y_1, \dots, Y_n) = \underbrace{P(X_1, \dots, X_n)}_{\text{Prior}} \underbrace{P(Y_1, \dots, Y_n \mid X_1, \dots, X_n)}_{\text{Likelihood}}$$

Making observations



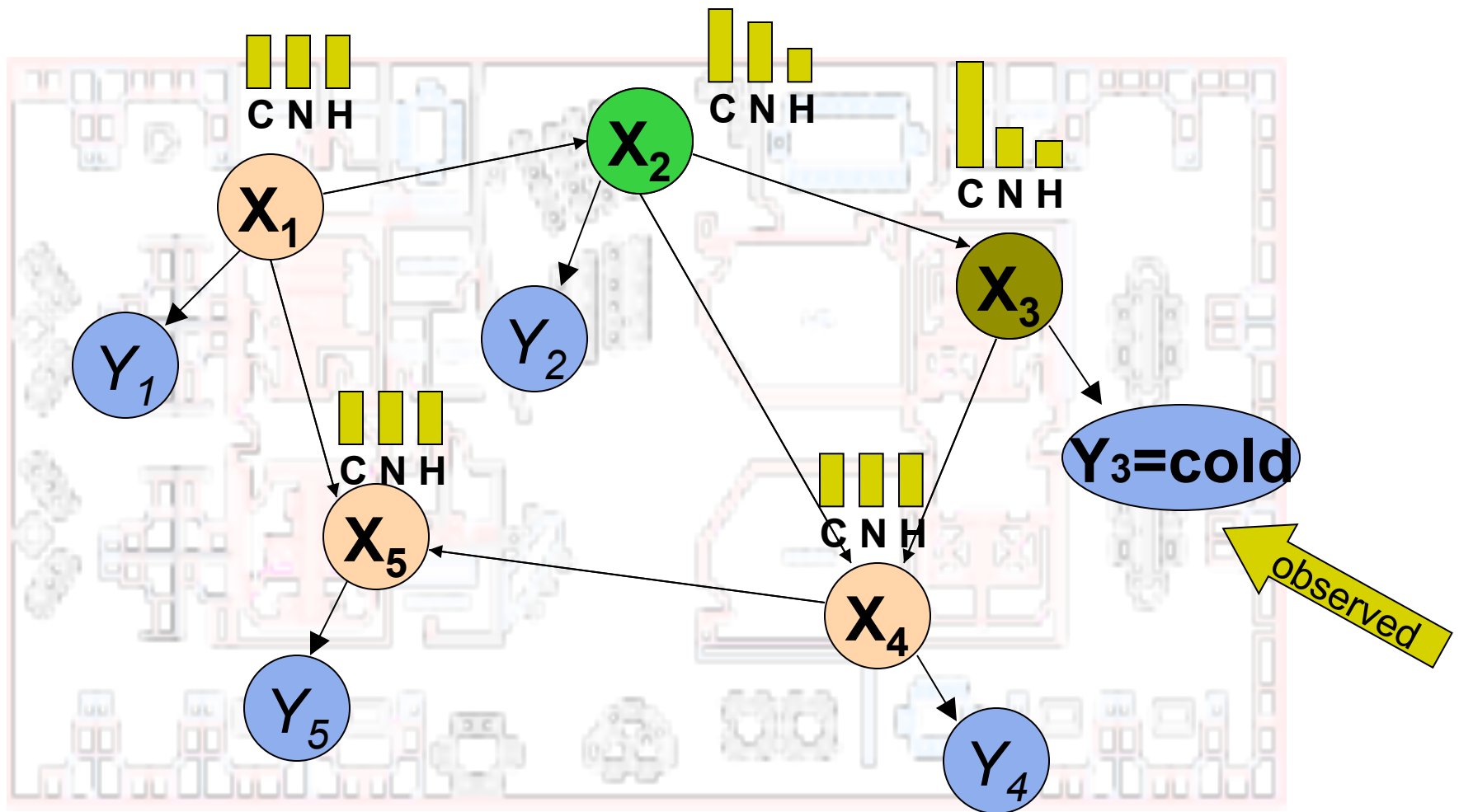
Less uncertain $\rightarrow \text{Reward}[P(\mathbf{X} | Y_1 = \text{hot})] = 0.2$

Making observations



$$\text{Reward}[P(\mathbf{X} | Y_3 = \text{hot})] = 0.4$$

A different outcome...



$$\text{Reward}[P(\mathbf{X} | Y_3 = \text{cold})] = 0.1$$

Reducing uncertainty

- Want to select observations that maximize reduction in uncertainty

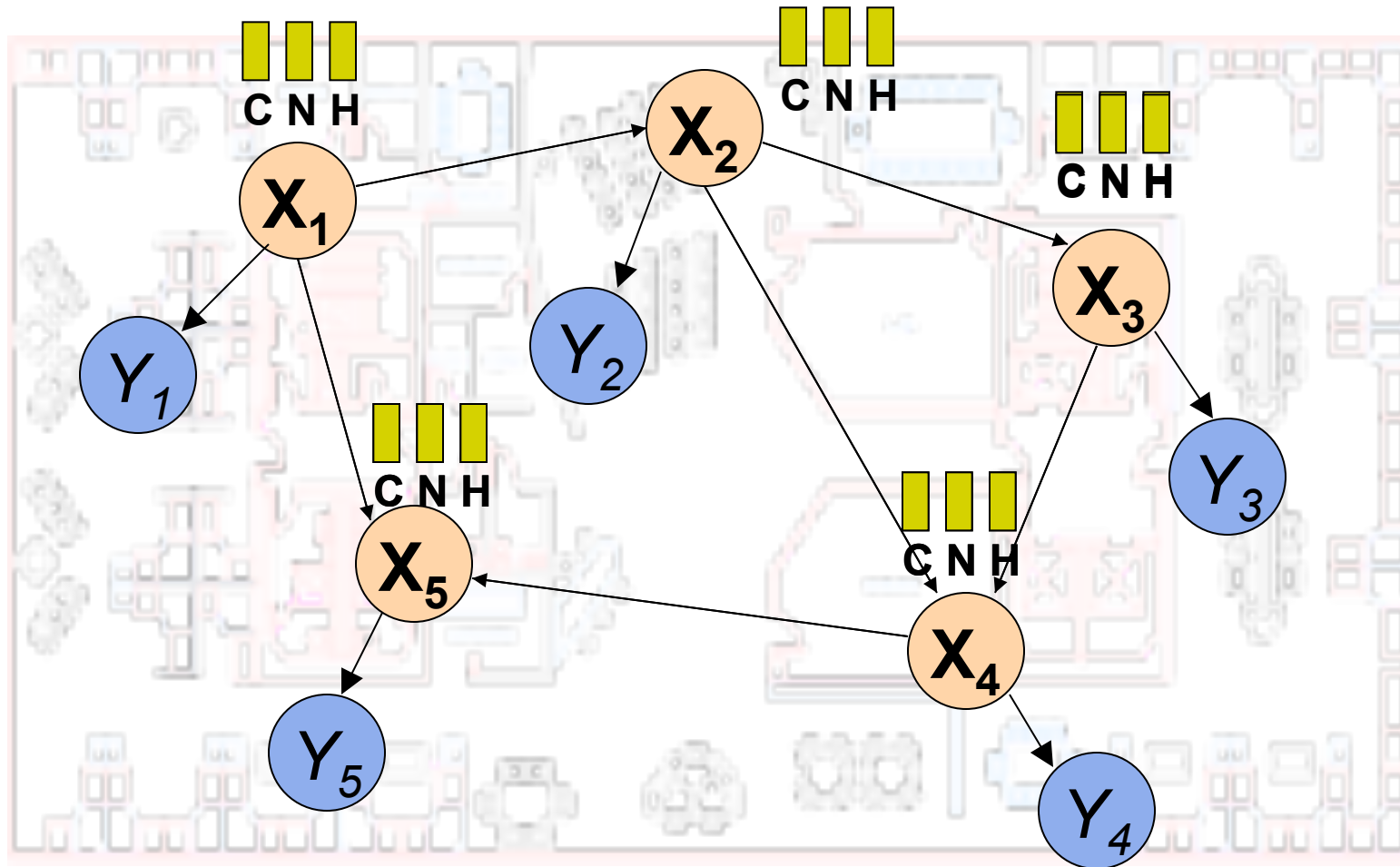
- Can quantify uncertainty using Shannon entropy:

$$H(X) = - \sum_x P(X = x) \log_2 P(X = x)$$

- For discrete variables $0 \leq H(X) \leq \log_2 |dom(X)|$
Spec. $P(X=x) = \frac{1}{n} \Rightarrow H(X) = -n \cdot \frac{1}{n} \log_2 \frac{1}{n} = \log_2 n$

- Thus, can use $\text{Reward}[P(\mathbf{X})] = -H(\mathbf{X}) = \sum_x P(\mathbf{x}) \log_2 P(\mathbf{x})$

Making observations



Prior entropy: $H(\mathbf{X}) \approx 4.2$

Posterior entropy

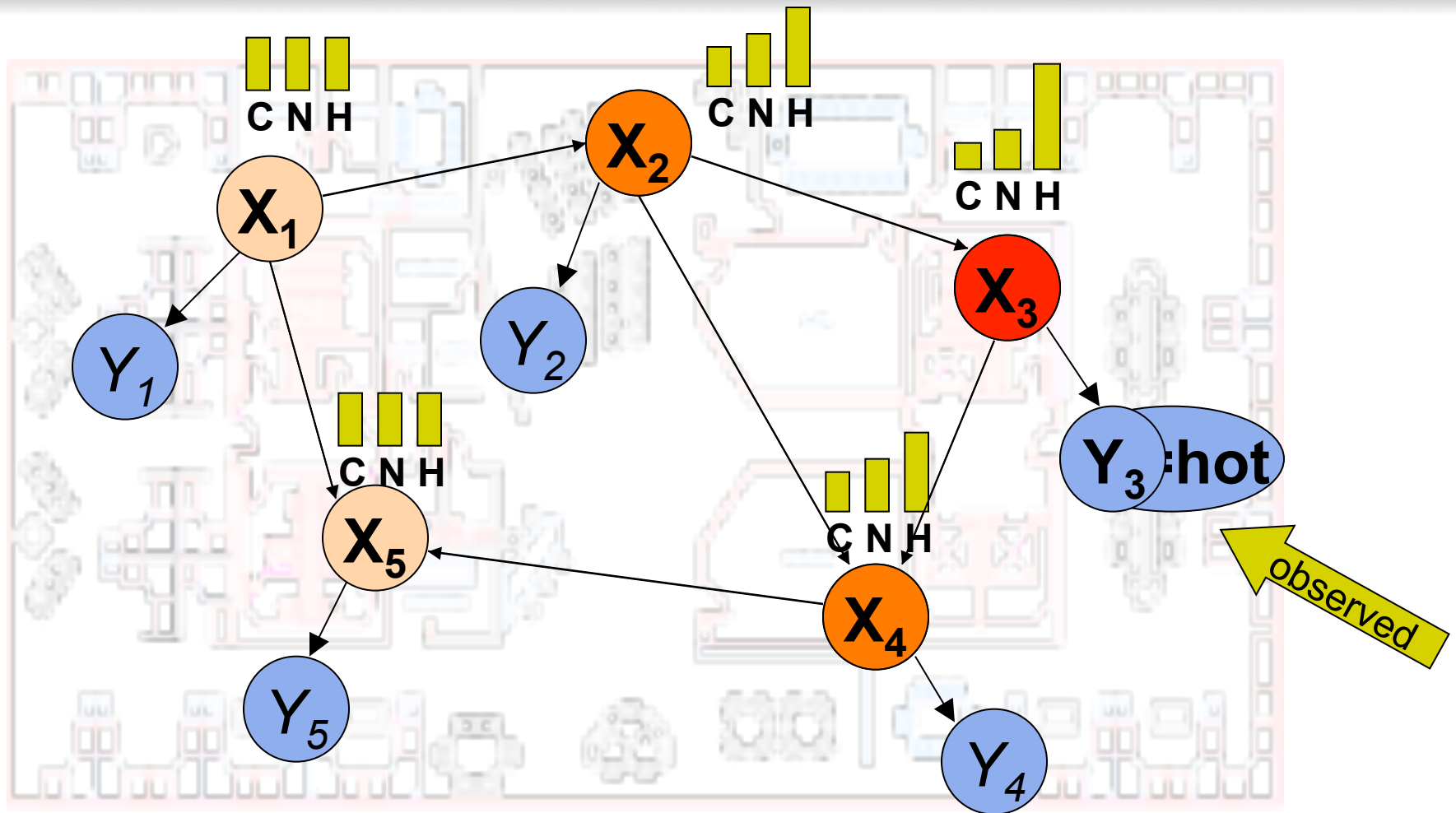
- Entropy before observations:

$$H(X) = - \sum_x P(X = x) \log_2 P(X = x)$$

- Entropy *after* observing $Y = y$:

$$H(X \mid Y = y) = - \sum_x P(X = x \mid Y = y) \log_2 P(X = x \mid Y = y)$$

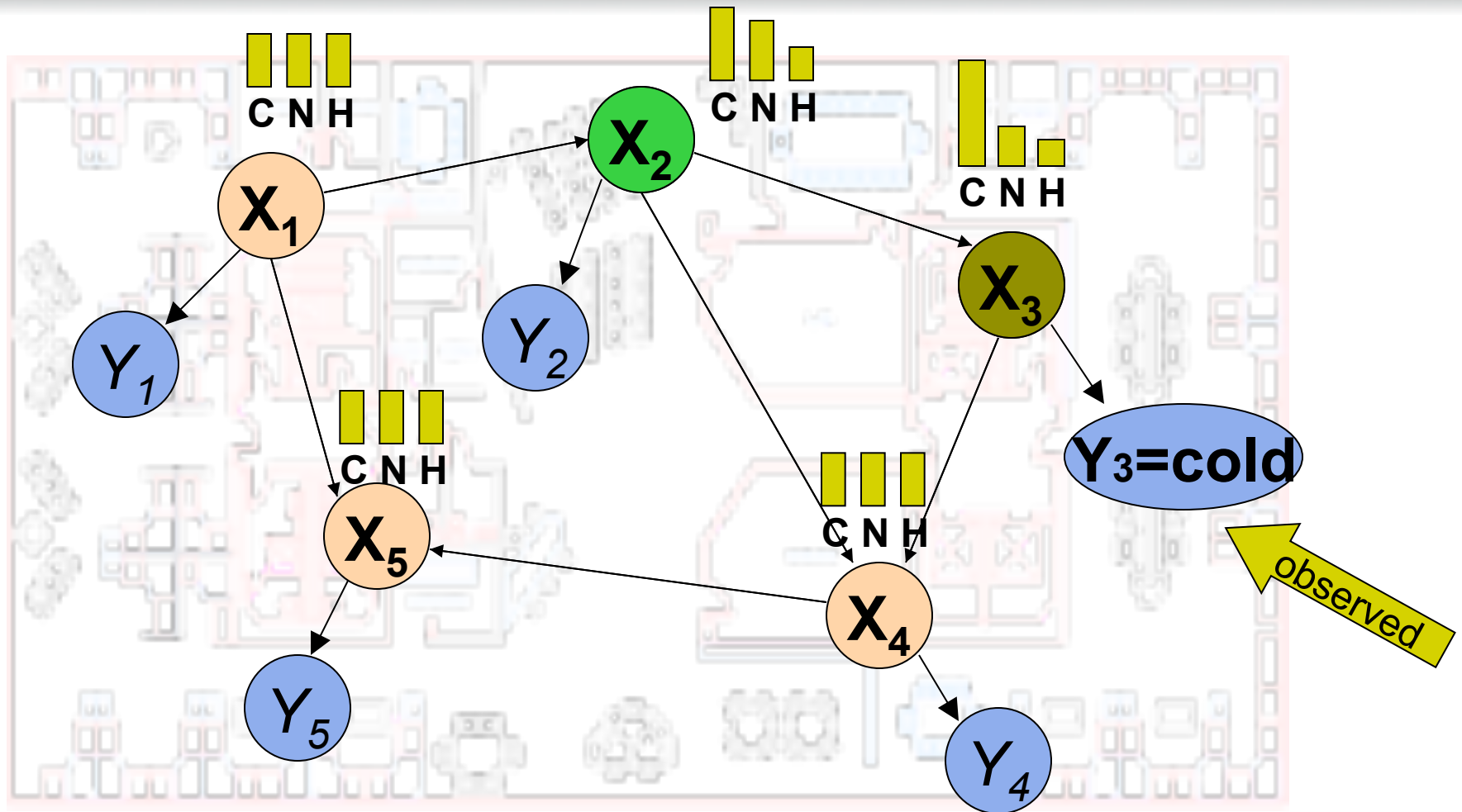
Making observations



Posterior entropy $H(\mathbf{X} \mid Y_3 = \text{hot}) \approx 2.7$

Reward: $H(\mathbf{X}) - H(\mathbf{X} \mid Y_3 = \text{hot}) \approx 1.5$

A different outcome...



Posterior entropy $H(\mathbf{X} \mid Y_3 = \text{cold}) \approx 3.2$

Reward: $H(\mathbf{X}) - H(\mathbf{X} \mid Y_3 = \text{cold}) \approx 1.0$

Information gain

- Entropy after observing $Y = y$:

$$H(X \mid Y = y) = - \sum_x P(X = x \mid Y = y) \log_2 P(X = x \mid Y = y)$$

- Don't know value of y before observing it!
- Conditional entropy:

$$H(X \mid Y) = \sum_y P(y) H(X \mid Y = y)$$

- Expected information gain (aka mutual information):

$$I(X; Y) = H(X) - H(X \mid Y)$$

Properties of entropy and infogain

$$\text{Prod. rule: } P(X, Y) = P(X) \cdot P(Y|X)$$

$$H(X, Y) = H(X) + H(Y|X)$$

$$H(X_1, \dots, X_n) = H(X_1) + H(X_2|X_1) + H(X_3|X_2, X_1) + \dots + H(X_n|X_1, \dots, X_{n-1})$$

$$H(X) \geq 0$$

$$\Rightarrow H(X) \geq H(X|Y) \quad \text{"information never hurts" (on average)}$$

$$I(X; Y) \geq 0$$

$$I(X; Y) = 0 \quad \text{iff } X \perp Y$$

$$\begin{aligned} I(X; Y) &= H(X) - \underbrace{H(X|Y)}_{H(X, Y) - H(Y)} = H(X) + H(Y) - \underbrace{H(X, Y)}_{H(Y, X)} \\ &= I(Y; X) \end{aligned}$$

Maximizing information gain

- Given: finite set V of locations

$$F(A) = \mathbb{I}(X; Y_A)$$

- Want: $A^* \subseteq V$ such that
$$A^* = \operatorname{argmax}_{|A| \leq k} F(A)$$

Typically NP-hard!

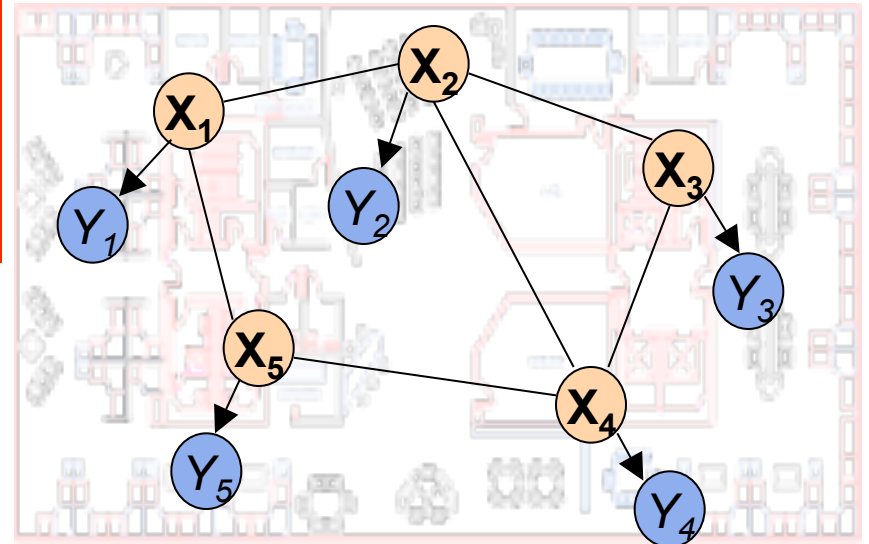
Greedy algorithm:

Start with $A = \{\}$

For $i = 1$ to k

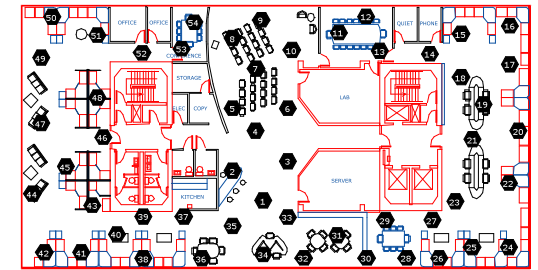
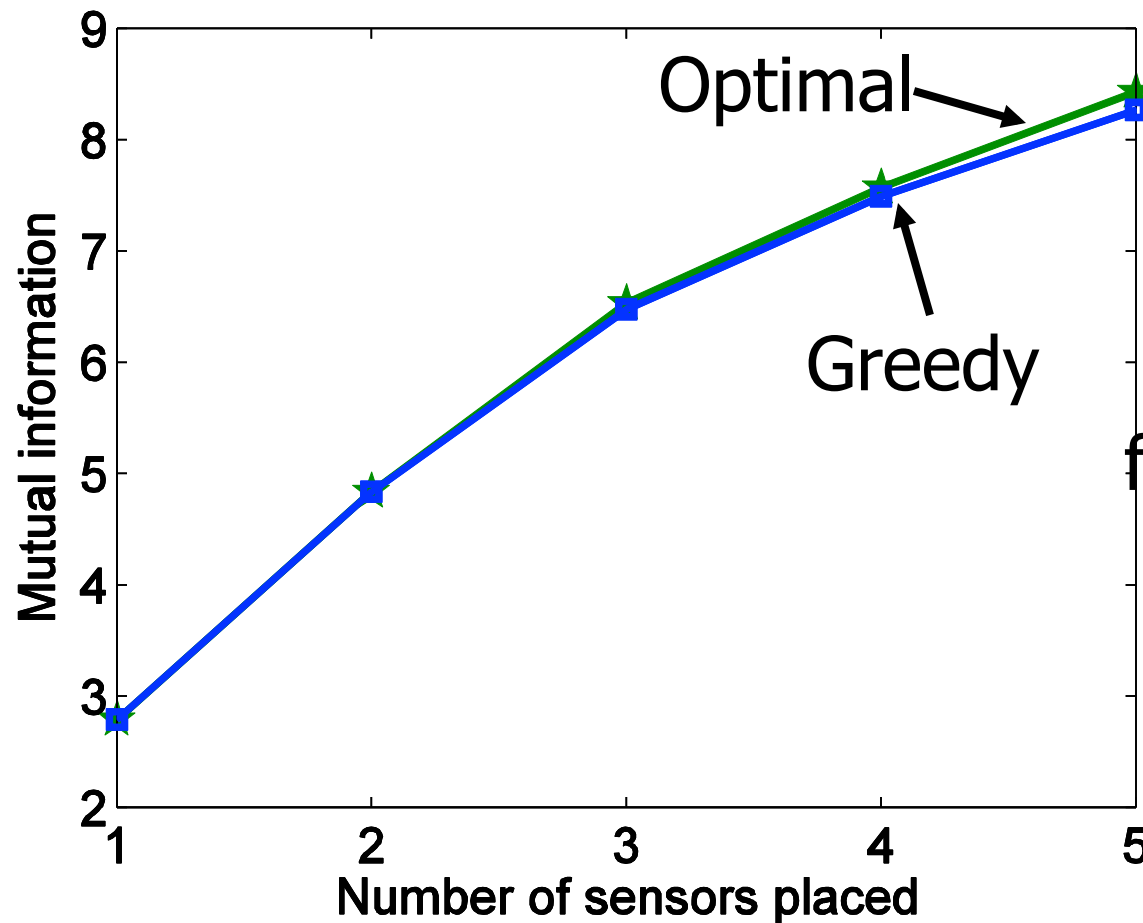
$s^* := \operatorname{argmax}_s F(A \cup \{s\})$

$A := A \cup \{s^*\}$



How well can this simple heuristic do?

Performance of greedy

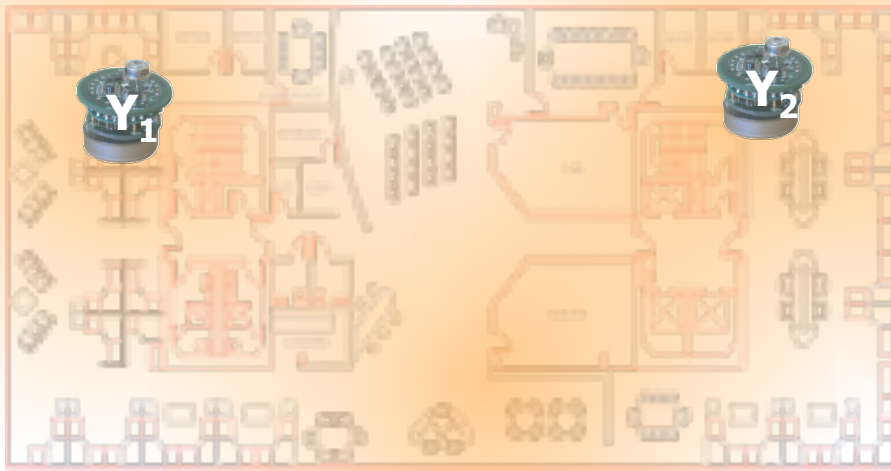


Temperature data
from sensor network

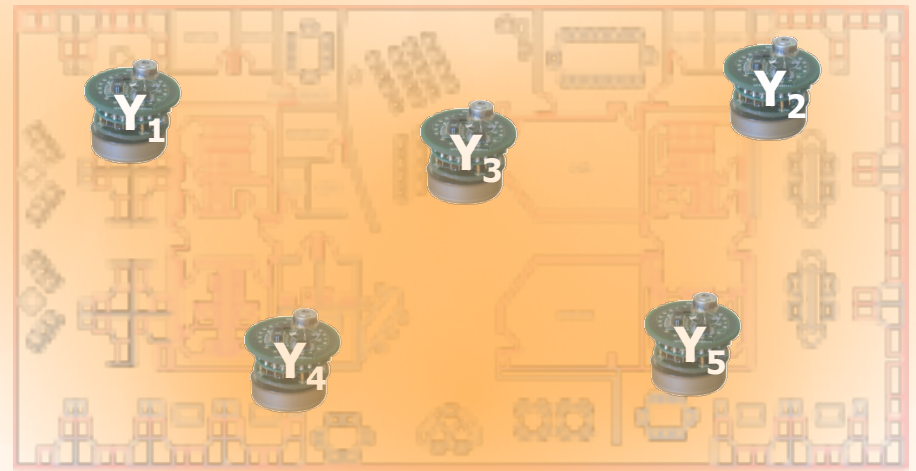
- Greedy empirically close to optimal. Why?

Key observation: Diminishing returns

Placement A = $\{Y_1, Y_2\}$



Placement B = $\{Y_1, \dots, Y_5\}$



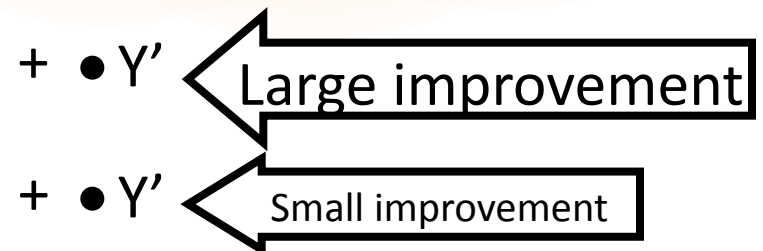
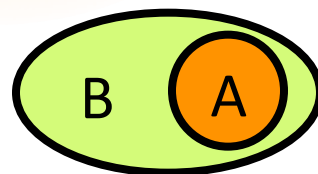
Adding Y' will help a lot!



Adding Y' doesn't help much

New sensor Y'

Submodularity:



$$\text{For } A \subseteq B, F(A \cup \{Y'\}) - F(A) \geq F(B \cup \{Y'\}) - F(B)$$

One reason submodularity is useful

Theorem [Nemhauser et al '78]

Greedy algorithm gives constant factor approximation

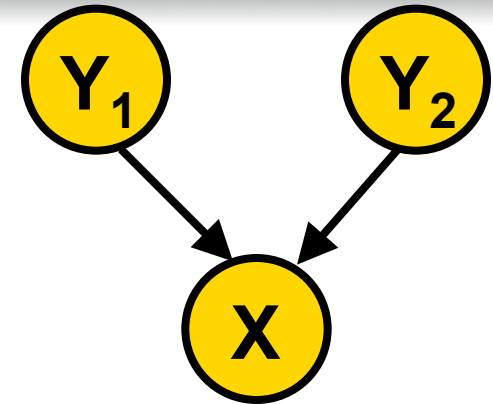
$$F(A_{\text{greedy}}) \geq (1-1/e) F(A_{\text{opt}})$$

- Greedy algorithm gives near-optimal solution!
- Is information gain submodular?

Non-submodularity of information gain

$Y_1, Y_2 \sim \text{Bernoulli}(0.5)$

$X = Y_1 \mathbf{XOR} Y_2$



Let $F(A) = I(Y_A; X) = H(X) - H(X|Y_A)$

$$X \sim \mathcal{B}(0.5) \quad H(X) = 1$$

$$X|Y_i = y \sim \mathcal{B}(0.5) \quad H(X|Y_i) = 1$$

$$X = Y_1 \mathbf{XOR} Y_2 \quad H(X|Y_1, Y_2) = 0$$

$$F(\emptyset) = H(X) - H(X) = 0$$

$$F(\{Y_1\}) = H(X) - H(X|Y_1) = 0$$

— " — $Y_2 = 0$

$$F(\{Y_1, Y_2\}) = \overset{H(X)}{H(X|Y_1, Y_2)} = 1$$



Example: Submodularity of info-gain

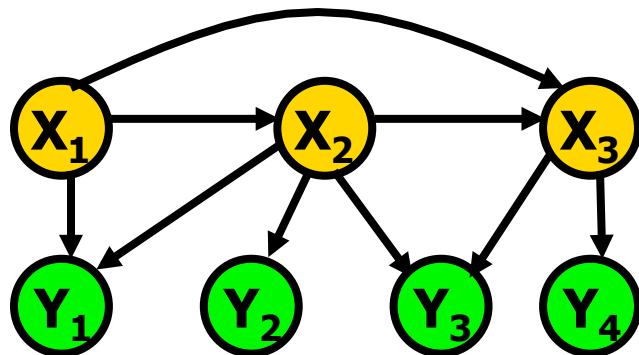
$Y_1, \dots, Y_m, X_1, \dots, X_n$ discrete RVs

$$F(A) = I(X; X_A) = H(Y) - H(Y | X_A)$$

- However, NOT always submodular

Theorem

If Y_i are all conditionally independent given X ,
then $F(A)$ is submodular!



Hence, greedy algorithm works!

In fact, NO algorithm can do better
than $(1-1/e)$ approximation!

Case study: Building a Sensing Chair

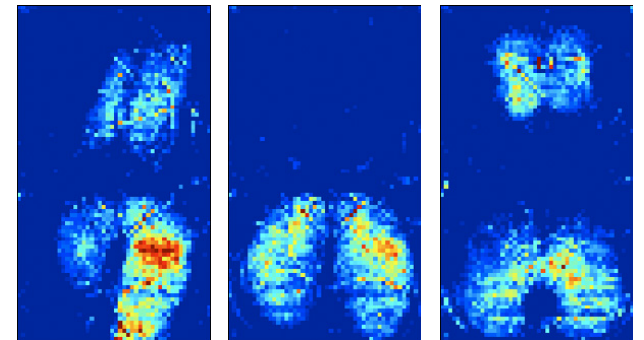
- Activity recognition in assistive technologies
- Seating pressure as user interface



Equipped with
1 sensor per cm²!

Costs \$6,000!

Can we get similar
accuracy with fewer,
cheaper sensors?



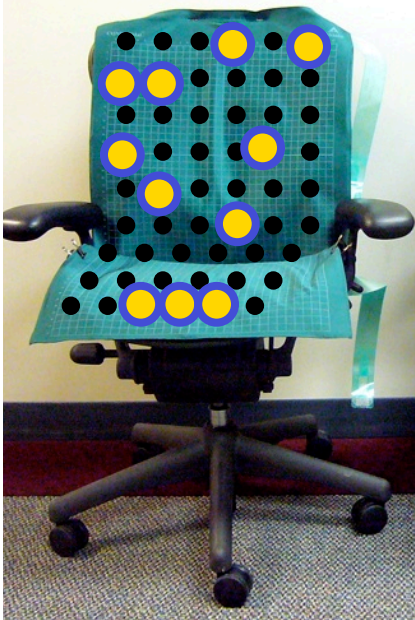
Lean Lean Slouch
left forward
**82% accuracy on
10 postures!**

How to place sensors on a chair?

- Sensor readings at locations V as random variables
- Predict posture X using probabilistic model $P(Y,V)$
- Pick sensor locations $A^* \subseteq V$ to minimize entropy:

$$A^* = \operatorname{argmax}_{|\mathcal{A}| \leq k} I(X; \mathbf{Y}_A)$$

Possible locations V

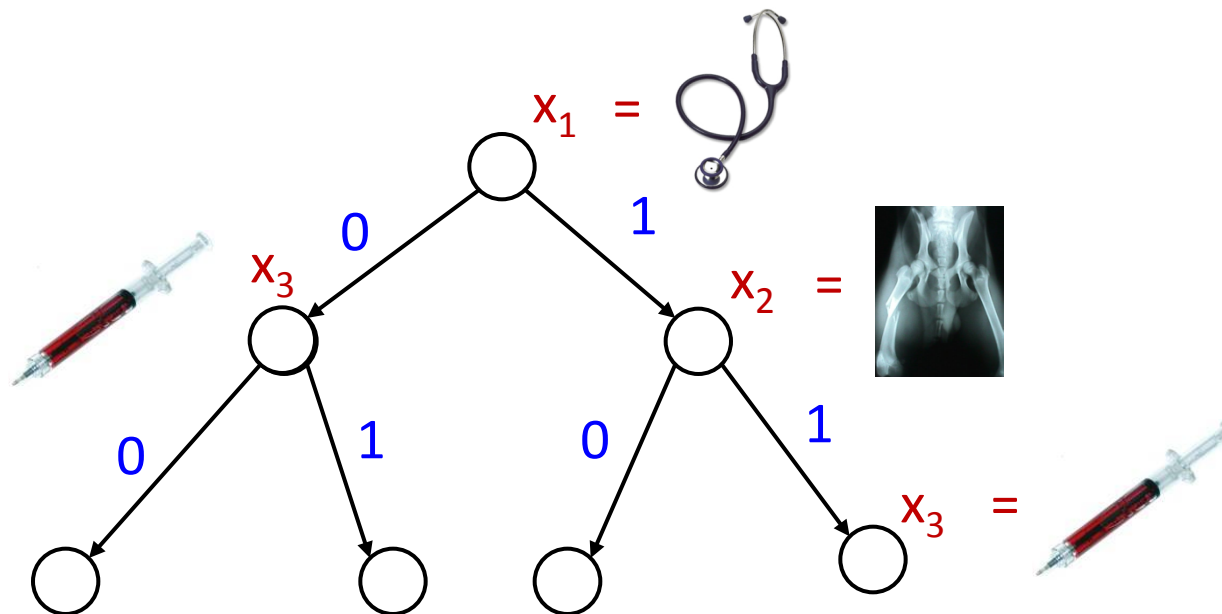


← Placed sensors, did a user study:

	Accuracy	Cost
Before	82%	\$6,000 ☹️
After		

Adaptive Optimization

- So far: Search for a most informative *set* of variables (e.g., sensor placement).
- In many applications, want to adaptively choose observations:



Interested in a *policy* (decision tree), not a *set*.

Adaptive greedy algorithm

- Expected benefit of adding test s after we've seen $Y_A = y_A$.

$$\Delta(s \mid \mathbf{y}_A) = H(\mathbf{X} \mid \mathbf{y}_A) - \sum_{y_s} P(y_s \mid \mathbf{y}_A) H(\mathbf{X} \mid \mathbf{y}_A, y_s)$$

Adaptive Greedy algorithm:

Start with $A = \emptyset$

For $i = 1:k$

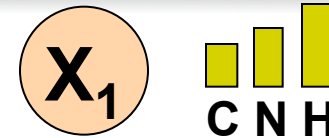
- Pick $s_k \in \arg \max_s \Delta(s \mid \mathbf{y}_A)$
- Observe $Y_{s_k} = y_{s_k}$
- Set $A \leftarrow A \cup \{s_k\}$

Gathering information for making decisions

- So far: Selecting variables which decrease the uncertainty the most
- Often, want to gather information to take the right action

Value of information

Should we raise a fire alert?



Actions \ Temp. X	<i>Fiery hot</i>	<i>normal/cold</i>
No alarm	-\$\$\$	0
Raise alarm	\$	-\$

Only have belief about temperature $P(X = \text{hot} \mid \text{obs})$

→ choose $a^* = \operatorname{argmax}_a \sum_x P(\mathbf{x} \mid \text{obs}) U(\mathbf{x}, a)$

Decision theoretic value of (perfect) information

$\text{Reward}[P(X \mid \text{obs})] = \text{MEU}(X \mid \text{obs}) = \max_a \sum_x P(\mathbf{x} \mid \text{obs}) U(\mathbf{x}, a)$

Value of information [Lindley '56, Howard '64]

For a set A of variables, its expected *value of information* is

$$F(A) = \sum_{\mathbf{y}_A} \underbrace{P(\mathbf{y}_A)}_{\substack{\text{Observations} \\ \text{made by sensors } \mathbf{A}}} \underbrace{\text{MEU}[\mathbf{X} \mid \mathbf{y}_A]}_{\substack{\text{Max. expected utility} \\ \text{when observing} \\ \mathbf{Y}_A = \mathbf{y}_A}}$$

Unfortunately, value of information is not submodular

Greedy algorithm can fail arbitrarily badly

Can do better with look-ahead

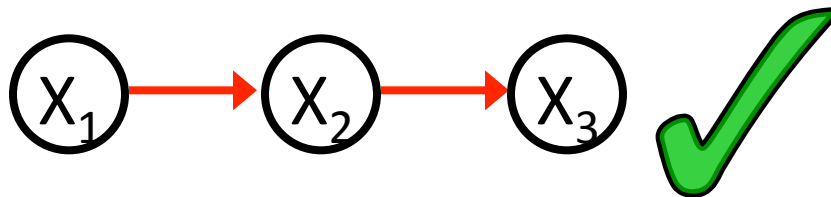
Maximizing value of information

[Krause, Guestrin '05]

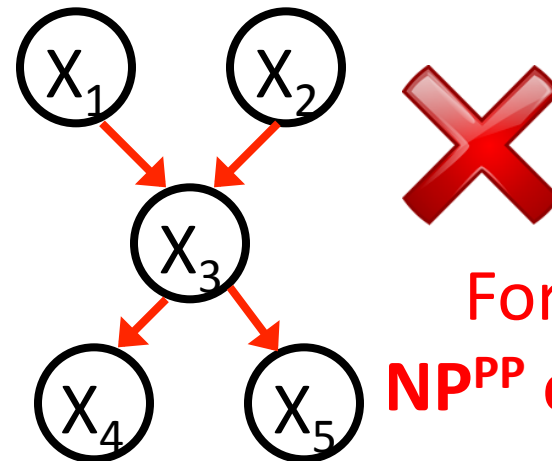
- Want to find a subset A^* of V , $|A^*| \leq k$ s.t.

$$A^* = \operatorname{argmax}_{|A| \leq k} F(A)$$

Theorem: Complexity of optimizing value of information



For chains (HMMs, etc.)
Optimally solvable in polytime 😊



For trees:
NP^{PP} complete 😞