

Hidden values

Expectation Maximization

Mixture Models

- Marked A4: now
- A5: Wednesday
- PP: Monday
- A6: out tonight, due next Wed.
- MT: 2 weeks from today
(study guide coming soon)

Why does bootstrap select $\approx 63\%$ for large 'N'?

$$p(\text{'i' selected at least once}) = 1 - p(\text{not selected 'N' times})$$

$$= 1 - \left(1 - \frac{1}{n}\right)^n$$

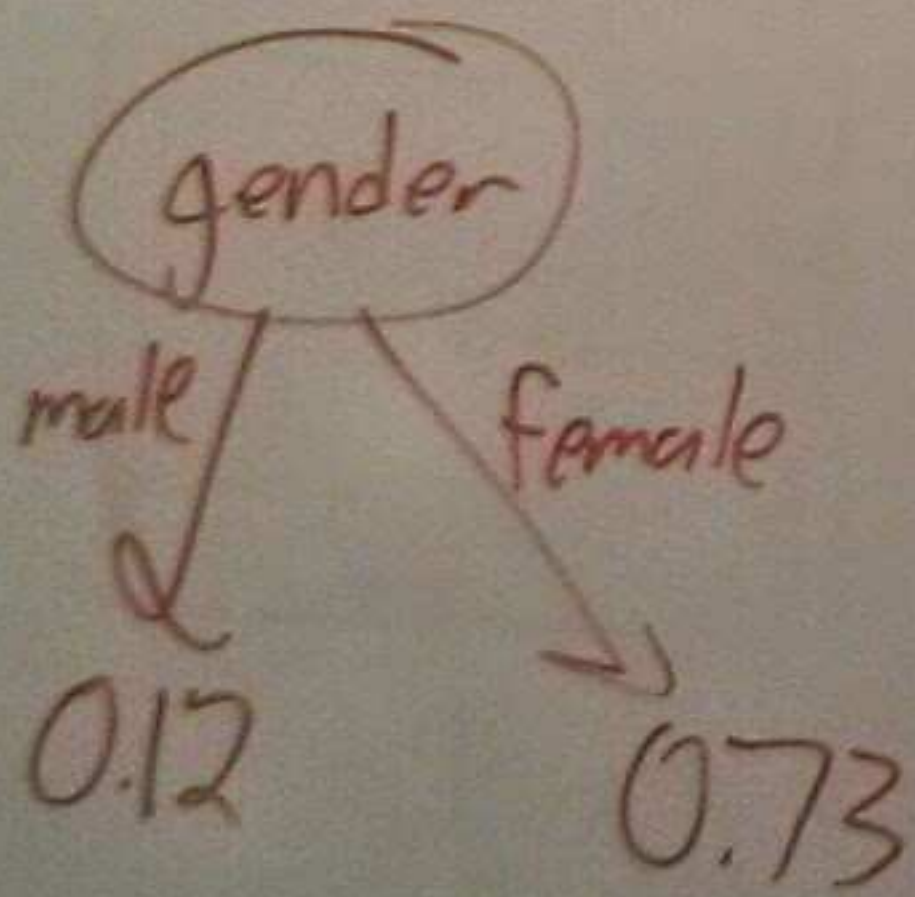
$$\approx 1 - \frac{1}{e}$$

$$\approx 63\%$$

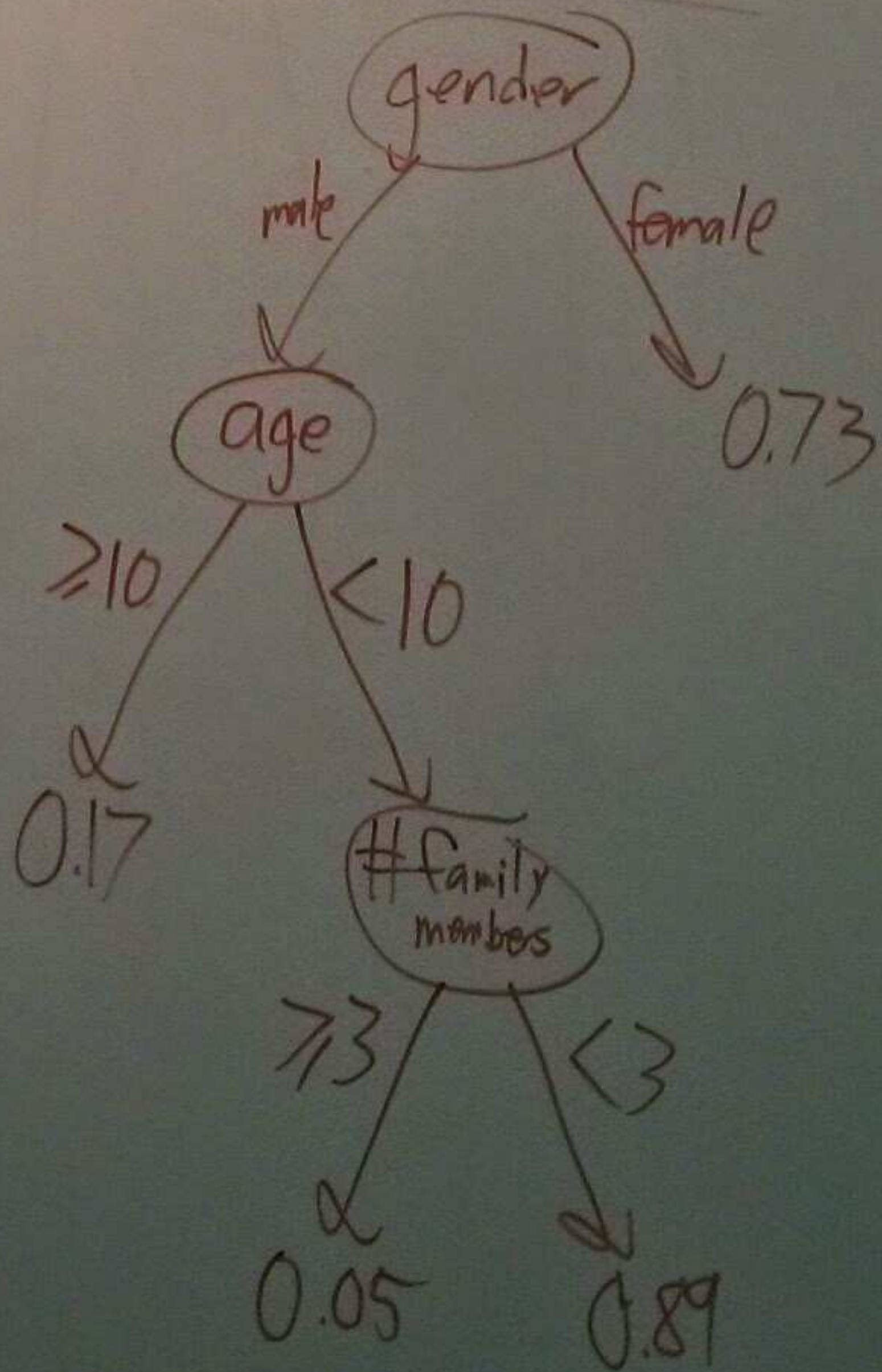
Example of decision tree

	Gender	Age	# Family members	
X =	Male	33	5	y =
	Female	10	1	
	⋮	⋮	⋮	
	⋮	⋮	⋮	
	⋮	⋮	⋮	

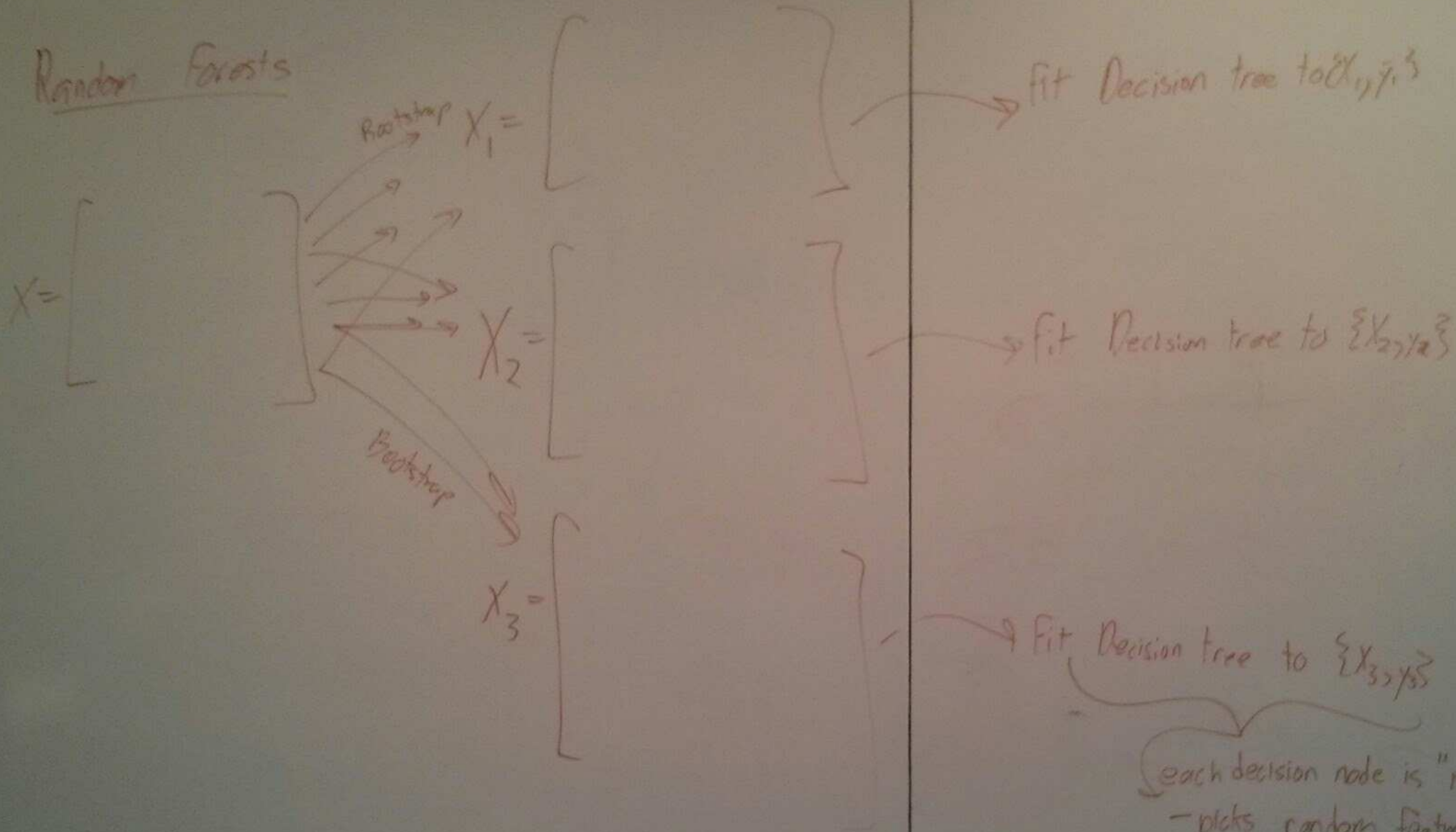
Decision Stump



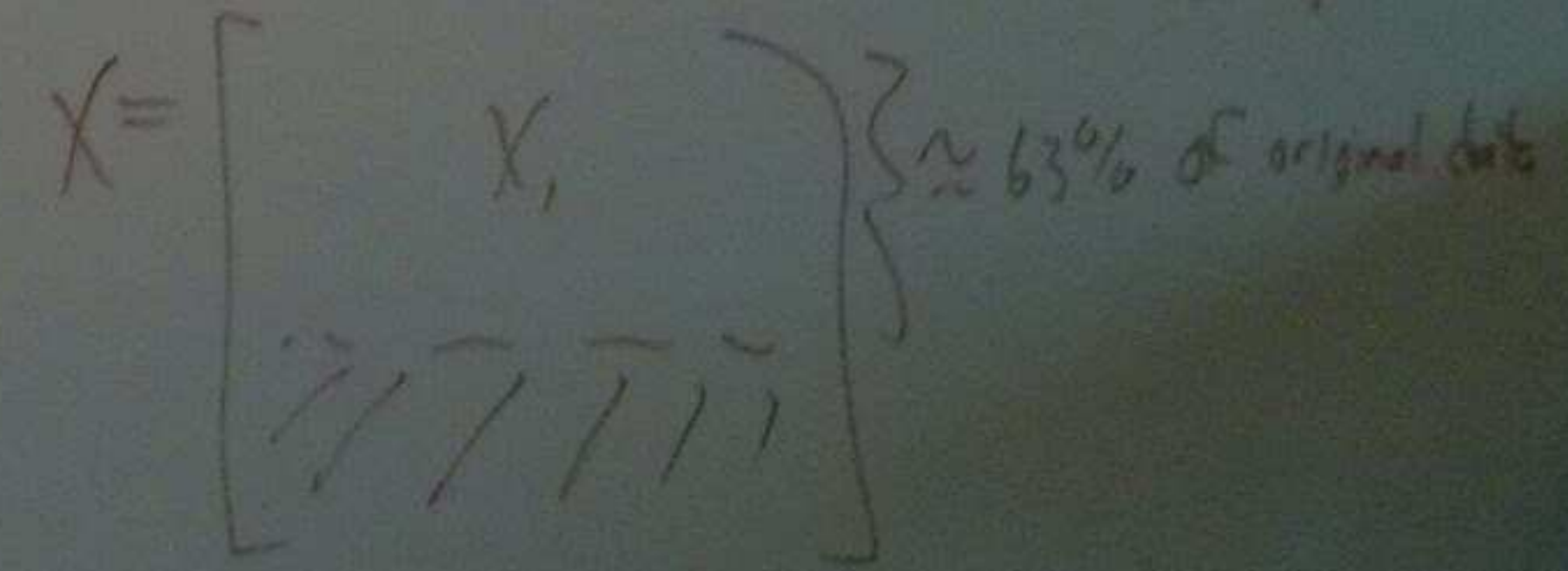
Decision tree



Random Forests



- each decision node is "random decision stump"
- picks random features
- use info gain
- you can prune with "left-out" samples



Hidden values

Expectation Maximization

Mixture Models

Hidden Values

- Learning when some values are
unobserved, missing, hidden, latent.

Example of decision tree

	Gender	Age	# family members	
X =	Male	33	5	y =
	Female	10	1	
	Female	?	2	
	Male	22	0	
				'lived'
				'died'
				?

Semi-Supervised Learning

Idea: getting labels is expensive,
getting unlabeled data is cheap.

Can we train on $\{X_L, Y_L\}$ and $\{X_U\}$

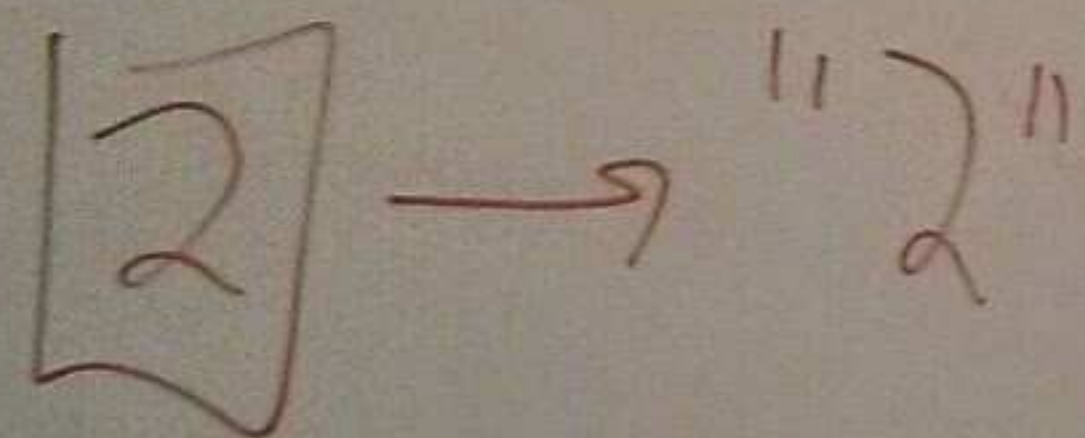
$X_L =$ [] $Y_L =$ []

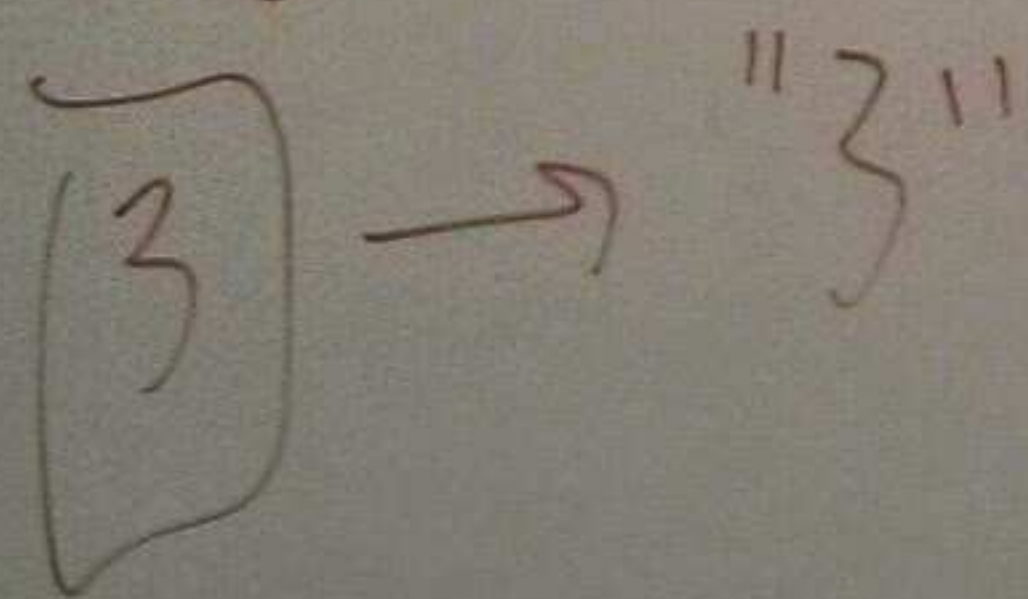
$X_U =$ [] $Y_U =$ []

"Missing at Random" (MAR)

- the fact that it is missing does not depend on the missing value.

- E.g. digit classification

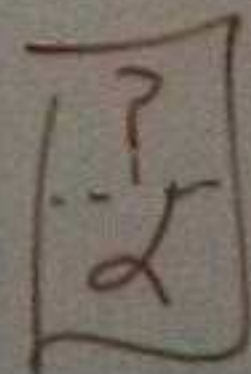




- missing random pixels: MAR

- hide the labels of all the "2" examples. (not MAR)

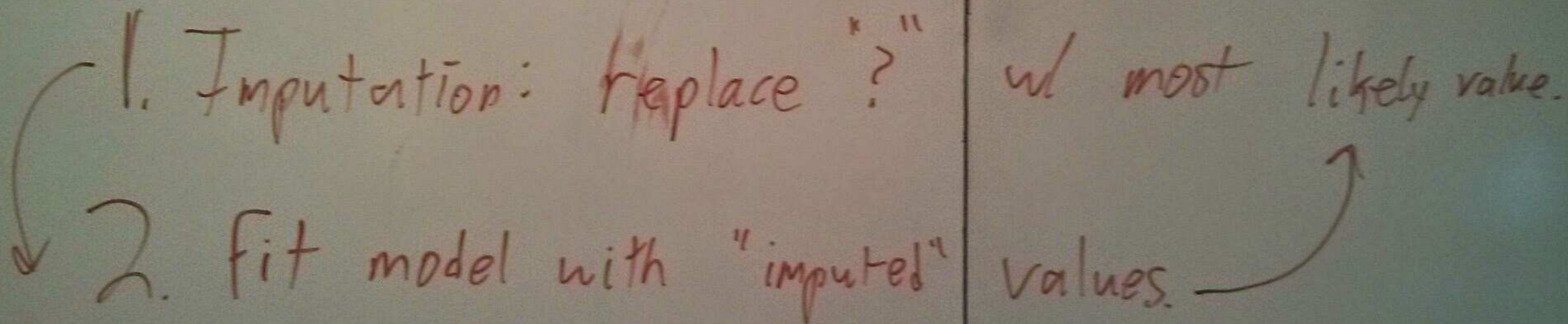
- hide the top half of every digit: MAR



- If not MAR, you need to model

WHY data is missing.

Approach #1

1. Imputation: replace "?" w/ most likely value.
 2. Fit model with "imputed" values.
- 

"hard-EM"

hidden values

Expectation Maximization

Mixture Models

Probabilistic Approach

Notation: X : observed variables

H : hidden variables

$$P(X) = \sum_h P(X, H=h)$$

(integral if H is a continuous r.v.)

E.g. SSL:

$$p(\bar{y}_L, X_L, X_U) = \prod_{i=1}^N p(y_i, x_i) \prod_{j=1}^T \left(\sum_{y_j} p(y_j, x_j) \right)$$

Problem:

Assume $-\log p(X, H)$ is "nice"
(closed-form, convex)

$$-\log(p(X)) = -\log\left(\sum_h p(X, H=h)\right)$$

not convex

"problem is Σ inside log"

$$\log(1 + \exp(w^T x))$$

$$\log(\exp(1) + \exp(w^T x)) \quad \text{"convex"}$$



Expectation Maximization

- local optimizer when $(\log p(X, H))$ is "nice" with parameters θ .

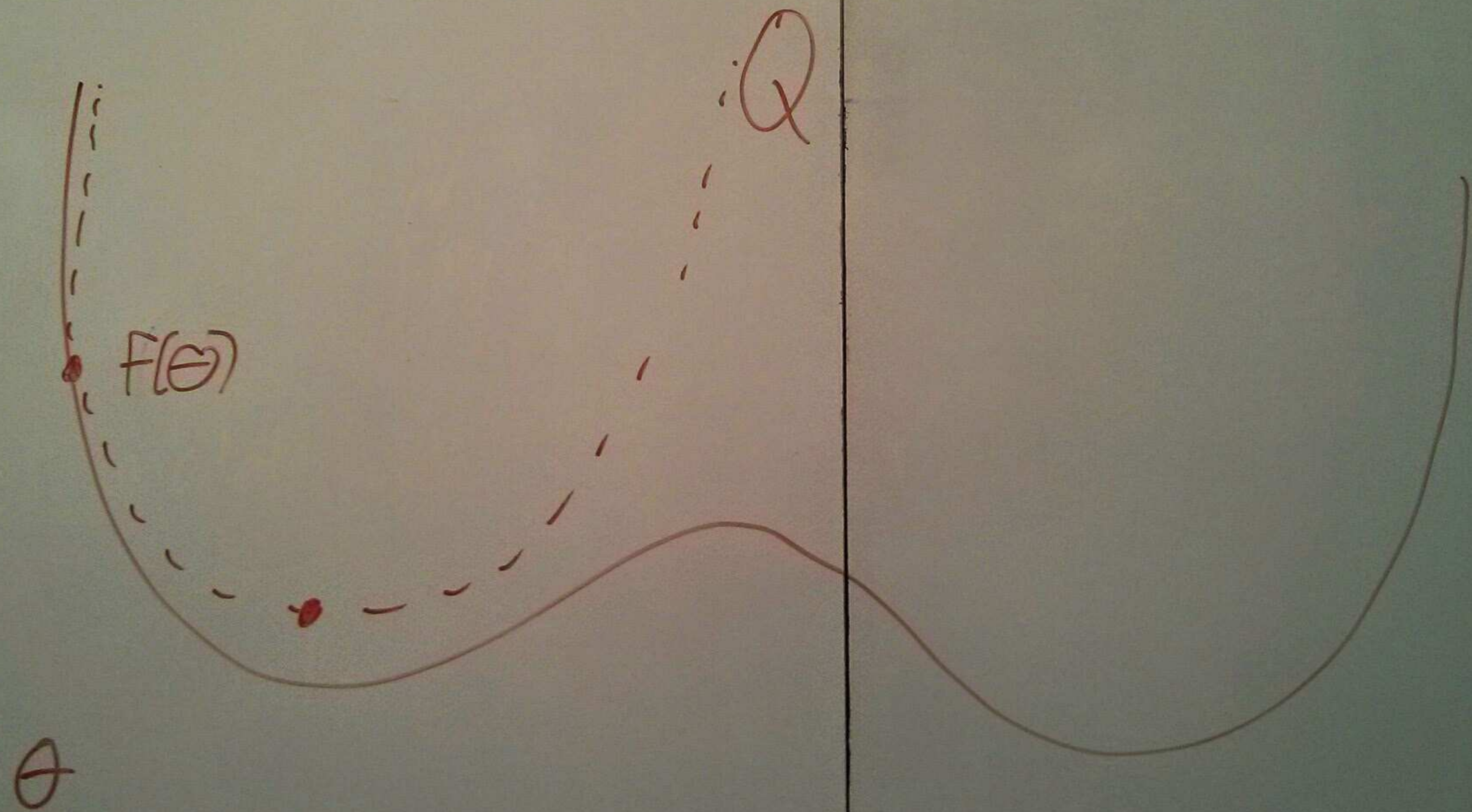
Problem: $\max_{\theta} p(X|\theta)$

Iterations: θ^t

"E"-step: Define $Q(\theta|\theta^t) = E_{H|X, \theta^t} [\log p(X, H|\theta)]$

"M"-step: $\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} Q(\theta|\theta^t) = \sum_h \underbrace{p(h|X, \theta^t)}_{\alpha_h^t} \underbrace{\log p(X, h|\theta)}_{\text{"nice"}}$

$$E_x[f(x)] = \sum_x p(x) f(x)$$
$$E_{x|y}[f(x)] = \sum_x p(x|y) f(x)$$



Maximization Models

Theoremi

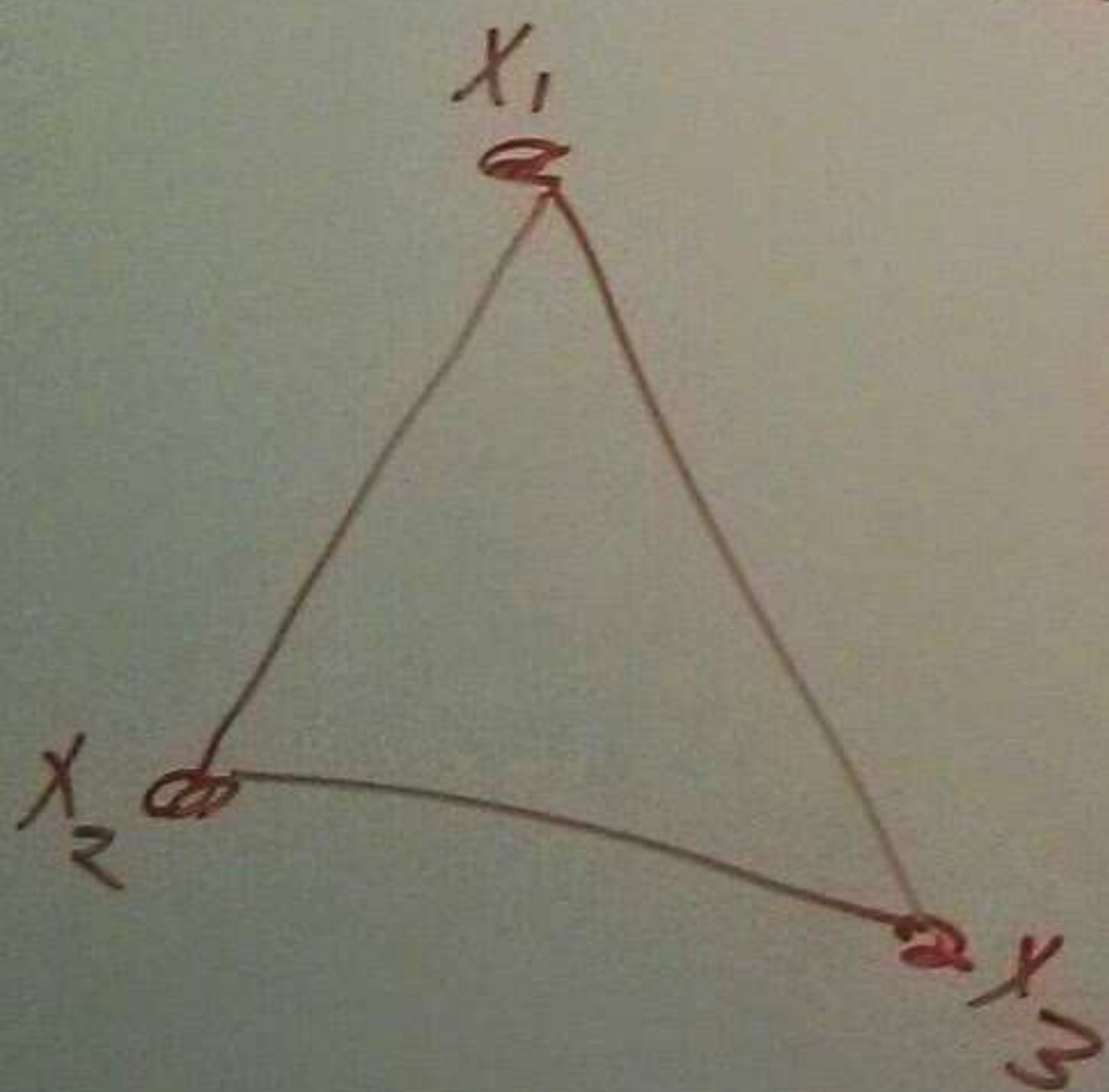
$$\log p(X | \theta^{t+1}) - \log p(X | \theta^t)$$

$$\geq Q(\theta^{t+1} | \theta^t) - Q(\theta^t | \theta^t)$$

"Convex Combination"

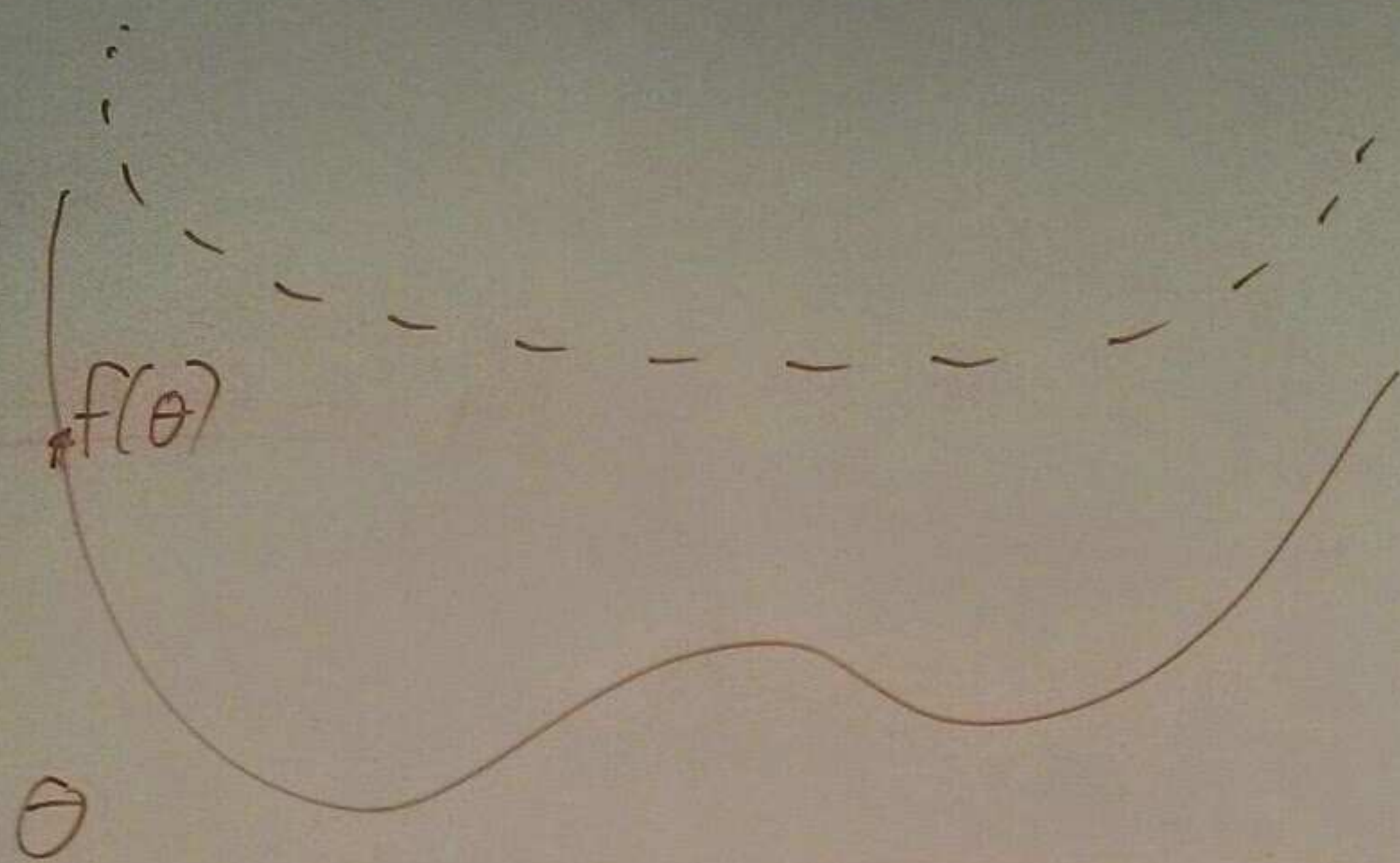
$$\alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \dots + \alpha_n x_n$$

$$\alpha_i \geq 0 \quad \sum \alpha_i = 1$$



If f is convex,

$$f\left(\sum_i \alpha_i x_i\right) \leq \sum_i \alpha_i f(x_i)$$



Problem: $\max_{\theta} p(X|\theta)$

Iterations: θ^t

"E"-step: Define $Q(\theta|\theta^t) = E_{H|X, \theta^t} [\log p(X, H|\theta)]$

"M"-step: $\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} Q(\theta|\theta^t) = \sum_h \underbrace{p(h|X, \theta^t)}_{\alpha_h^t} \underbrace{\log p(X, h|\theta)}_{\text{"nice"}}$

$$-\log p(X|\theta) = -\log \sum_h p(X, h|\theta)$$

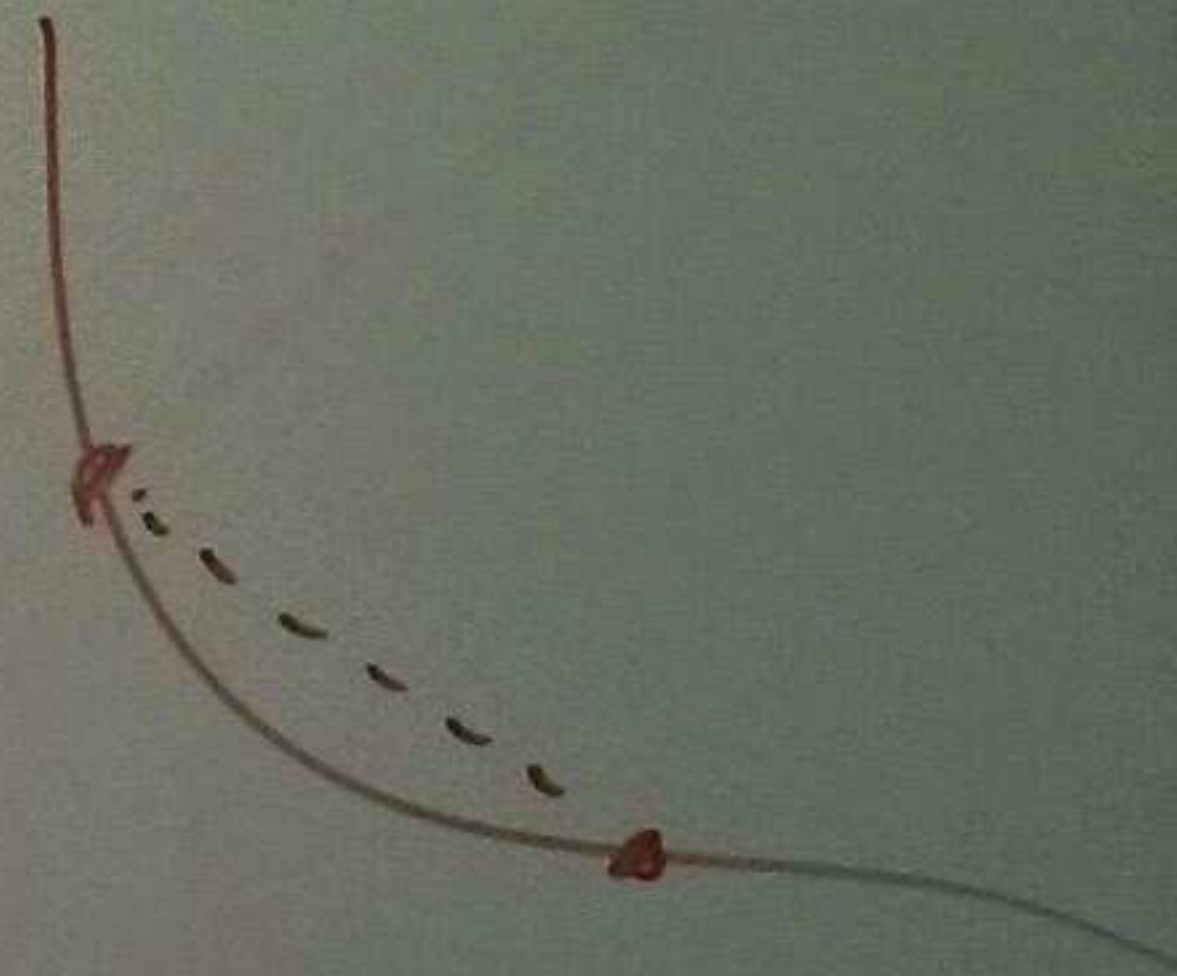
$$= -\log \left(\sum_h \alpha_h \frac{p(X, h|\theta)}{\alpha_h} \right)$$

$$\leq -\sum_h \alpha_h \log \frac{p(X, h|\theta)}{\alpha_h}$$

$$= -\sum_h \alpha_h \log p(X, h|\theta) + \sum_h \alpha_h \log \alpha_h$$

$$= -Q(\theta|\theta^t) + \text{const.}$$

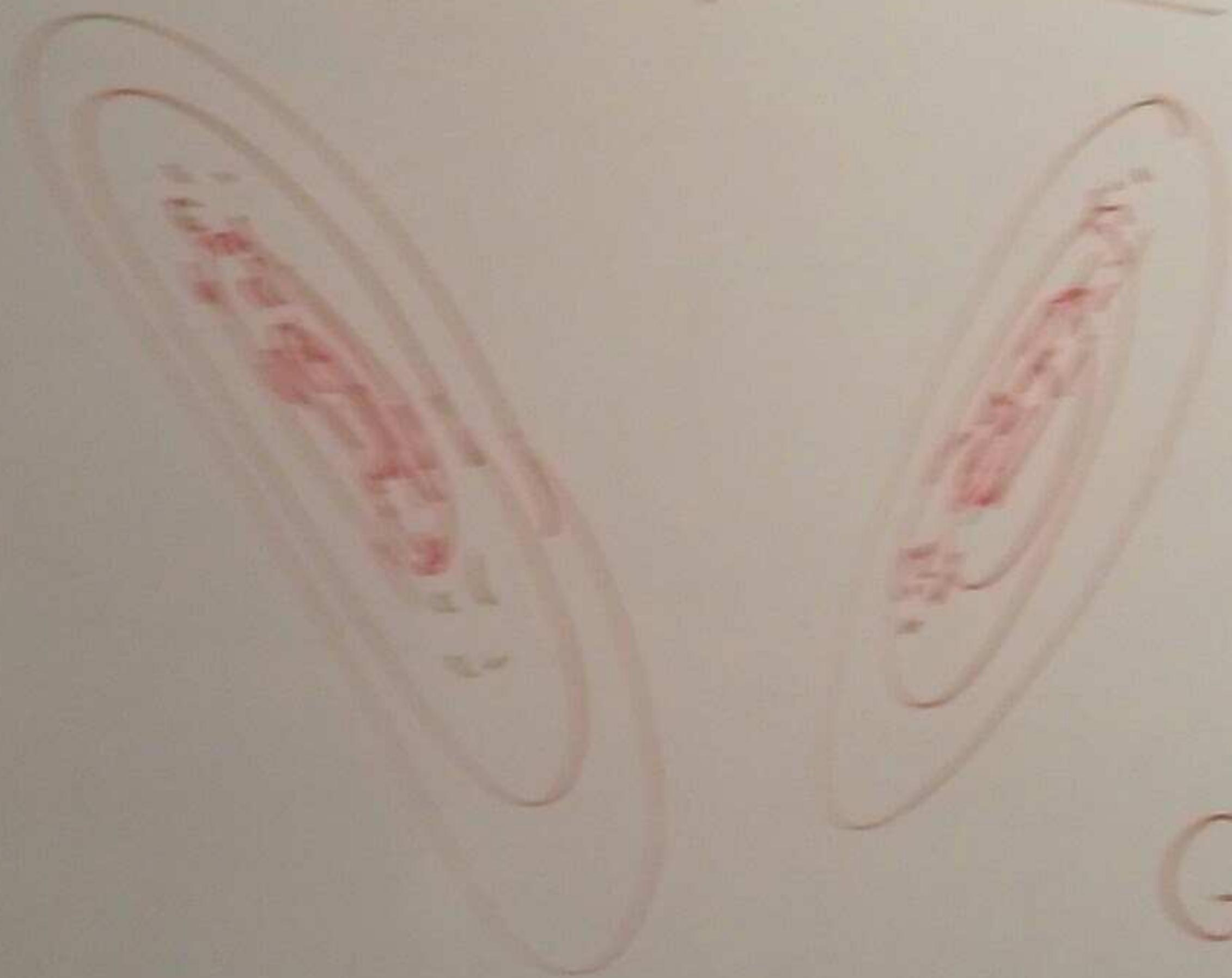
$$\underbrace{\hspace{10em}}_{-H(\vec{\alpha})}$$



Mixture Models

Mixture Models

Recall Fitting Gaussians:



GDA:

Want to fit probabilistic model:

- don't have labels

- but it would be easy if we did

- motivates "mixture" models

Let $z_i \in \{1, 2, \dots, K\}$

$$p_k(x_i | \theta) = p(x_i | z_i = k, \theta)$$

"Gaussian"

z_i is a "latent" variable.

$$p(x_i | \theta) = \sum_{k=1}^K p(x_i, z_i = k | \theta)$$
$$= \sum_{k=1}^K \underbrace{p(x_i | z_i = k, \theta)}_{\text{Gaussian for cluster 'k'}} \underbrace{p(z_i = k | \theta)}_{\text{probability that data point belongs to mixture}}$$

K-means: hard-EM
with Gauss
with $\sum_k = I$