

# Today

- 1. Admin, KNN
- 2. Probability
- 3. Naive Bayes

1. [www.cs.ubc.ca/~schmidtm/Courses](http://www.cs.ubc.ca/~schmidtm/Courses) /5

2. Piazza → sign up ASAP

3. Assignment 1 **DUE WEDNESDAY**  
(refresher, workload)

## 4. Alternative Courses

<u>CPSC</u>	<u>ECE</u>	<u>MATH</u>	<u>Stat</u>
322	592	302/303	305
340		318	447B
		418/419	460
			461
			547C

s/540

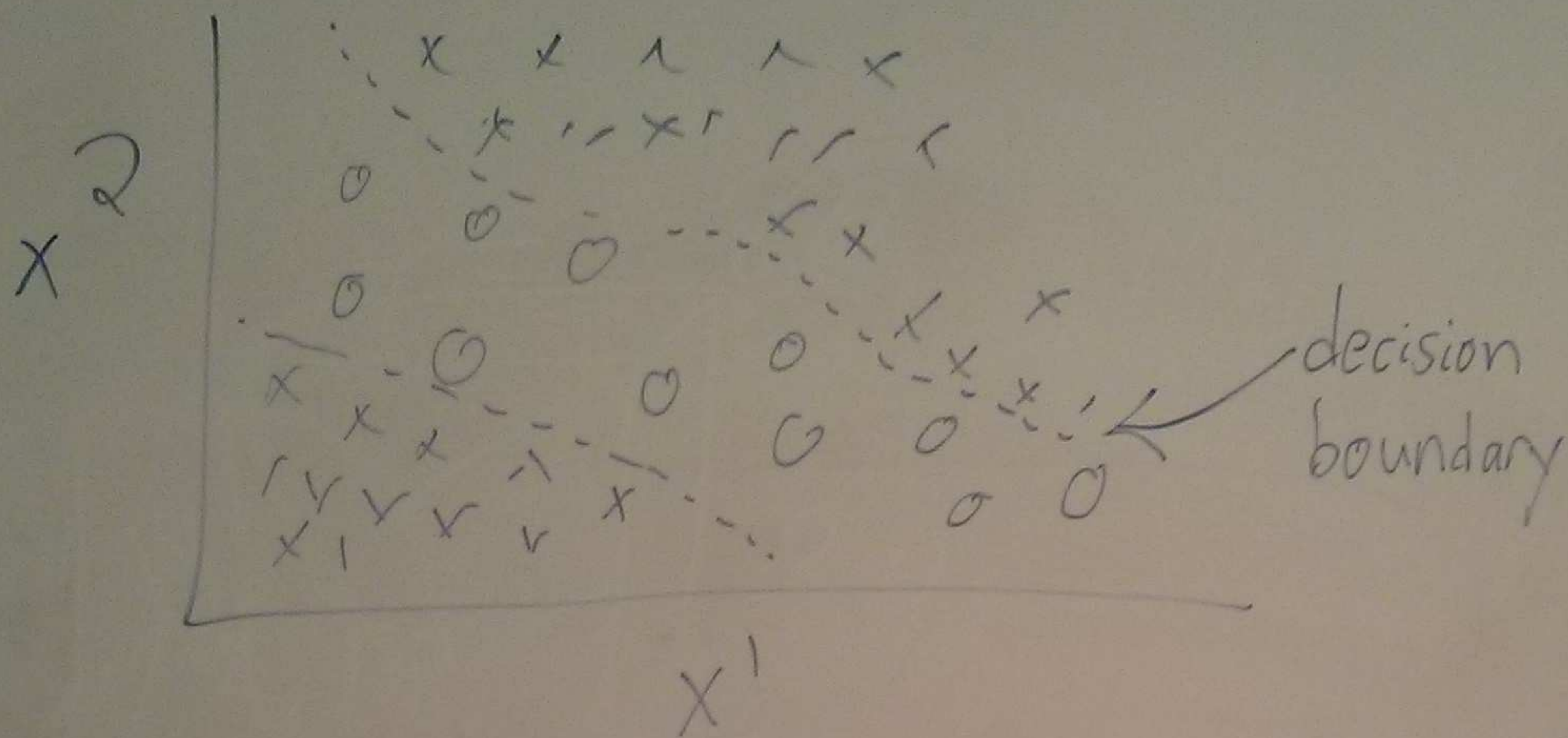
## Supervised Learning

Input:  $X, y$

Output: function  $f(x_i) \mapsto y_i$

Evaluation: compare  $f(x_i) = \hat{y}_i$  to  $y_i$  on new  $(x_i, y_i)$

Notation:  $x_i$  is training example 'i' (vector)  
 $x^i$  is variable 'i' (scalar)



Location based on features  $x^1, x^2$   
 +/0 based on class label  $y_i = +1$  or  $y_i = -1$

KNN: 'closest'

- Euclidean distance  $d(w, v) = \sqrt{\sum_{i=1}^D (w^i - v^i)^2} = \|w - v\|$   
 - others possible.

1-NN: 'cost'

- checking  $\|x_i - x_j\|$  is  $O(D)$

- do this  $N$  times for each of  $T$  test points:  $O(DNT)$

k-NN: 'cost'

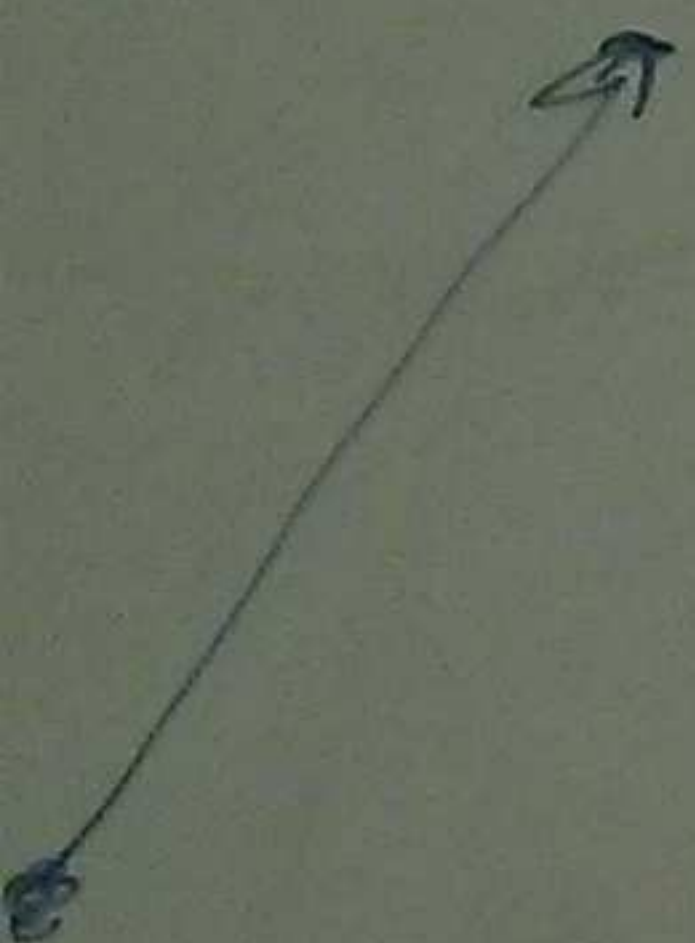
- sorting  $O(N \log N)$

- selecting  $O(kN)$

Possible course Project: Improve KNN

- better distance (learn it)

- faster classification



or  $y_i = -1$

Q: Why can we learn?

$$(w^i - v^i)^2 = \|w - v\|$$

# Today

1. Admin, KNN
2. Probability
3. Naive Bayes
4. Maximum Likelihood

# Probability

$$(x_i, y_i) \sim D$$

$$0 \leq P(A) \leq 1$$

"definitely not"                      "definitely"

$$P(\bar{A}) = 1 - P(A)$$

We'll use  $p(x) \triangleq p(X=x)$

$$\sum_{x \in X} p(x) = 1$$

$$p(A, B) = p(A \cap B)$$

$$p(A \cup B) = p(A) + p(B) - p(A, B)$$

Cond Probs

$$P(A|B) = \frac{P(A, B)}{P(B)} \quad (\text{for } P(B) > 0)$$

'Product rule'  $P(A, B) = P(A|B)P(B)$

Bayes rule  $= P(B|A)P(A)$

$$P(A|B, C) = \frac{P(B|A, C)P(A|C)}{P(B|C)}$$

$$\propto P(B|A)P(A)$$

∴ Prob of A given (B and C)



$$\sum_{i=1}^6 \frac{\prod_{j=1}^6 \Gamma(\alpha_j) p_j}{\Gamma(\sum_{j=1}^6 \alpha_j)} = \frac{\prod_{j=1}^6 \Gamma(\alpha_j)}{\Gamma(\sum_{j=1}^6 \alpha_j)}$$

$$P(3) = \frac{1}{6}$$

$$P(\neg 3) = 1 - P(3) = \frac{5}{6}$$

$$\sum_{i=1}^6 p_i = 1$$

$$\begin{aligned} P(1 \vee 2) &= P(1) + P(2) - P(1, 2) \\ &= \frac{1}{6} + \frac{1}{6} - 0 \\ &= \frac{1}{3} \end{aligned}$$

$$P(3, \text{even}) = 0$$

$$P(3, \text{odd}) = \frac{1}{6}$$

$$\begin{aligned} P(3 | \text{odd}) &= \frac{P(3, \text{odd})}{P(\text{odd})} \\ &= \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{3} \end{aligned}$$

Marg. Prob.

$$p(A) = \sum_b p(A, B=b)$$

$$p(\text{even}) = \sum_{i=1}^6 p(\text{even}, i)$$

$$= 0 + \frac{1}{6} + 0 + \frac{1}{6} + 0 + \frac{1}{6}$$

$$= \frac{1}{2}$$

We can always add extra conditioning



Conditional

Independence

$$X \perp Y | Z \iff p(X, Y | Z) = p(X | Z) p(Y | Z)$$

# Today

1. Admin, KNN
2. Probability
3. Naive Bayes
4. Maximum Likelihood

## Probabilistic Classifiers

### Discriminative

$$p(y_i | x_i)$$

- no model of  $x$
- "need lots of labelled data" but high accuracy
- e.g., logistic regression

### Generative

$$p(y_i, x_i)$$

- need to model  $x$
- "small amount of labelled data"
- e.g., Naive Bayes

### Training

$$p(x_i, y_i) = p$$

### Testing

$$p(y_i | x_i) =$$

$\alpha$

$y$   
 0  $p(y=0)$   
 1  $p(y=1)$

easy

Training

$$p(x_i | y_i) = p(x_i | y_i) p(y_i)$$

$$X = \begin{pmatrix} 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 1 & 1 & 0 \end{pmatrix} \quad y = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}$$

Testing

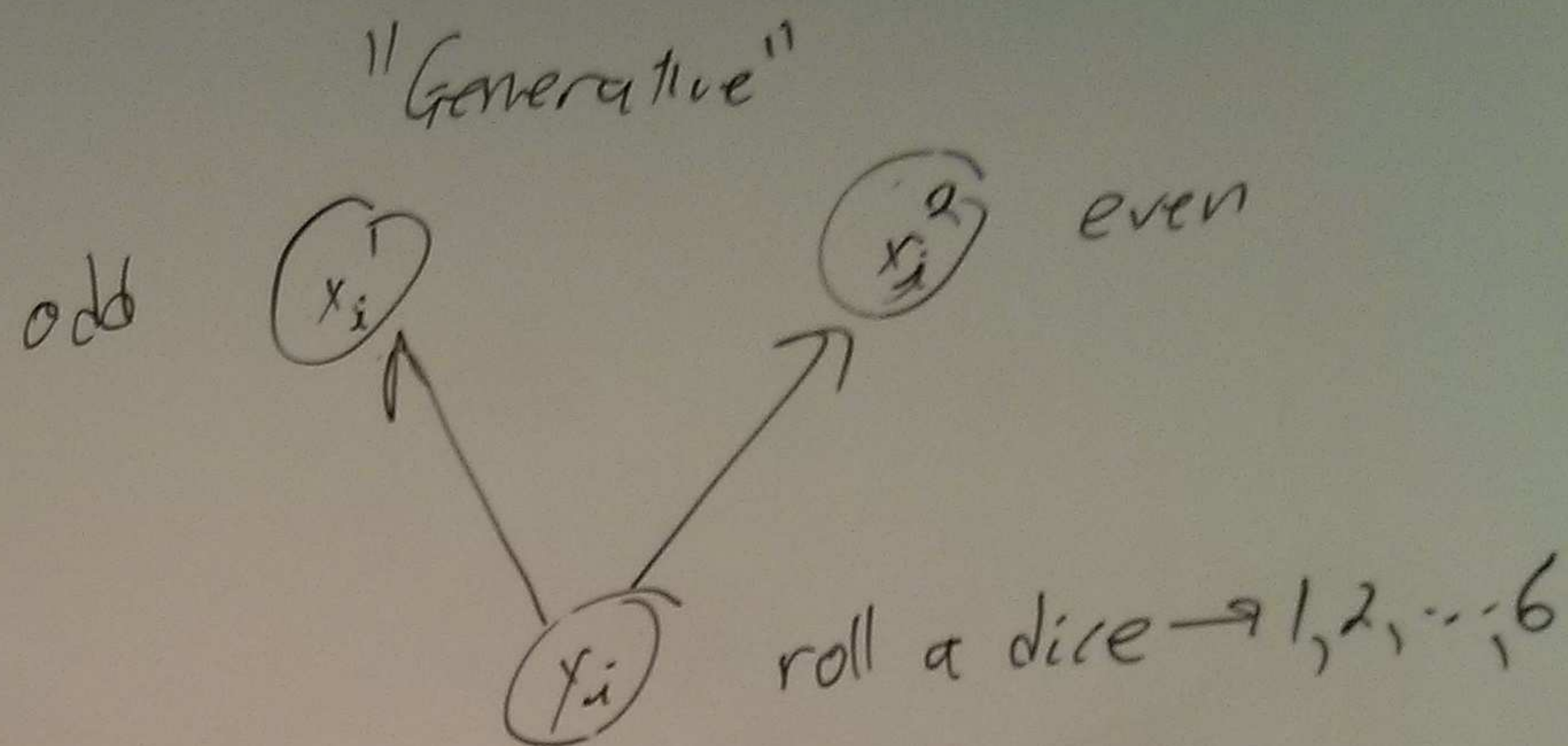
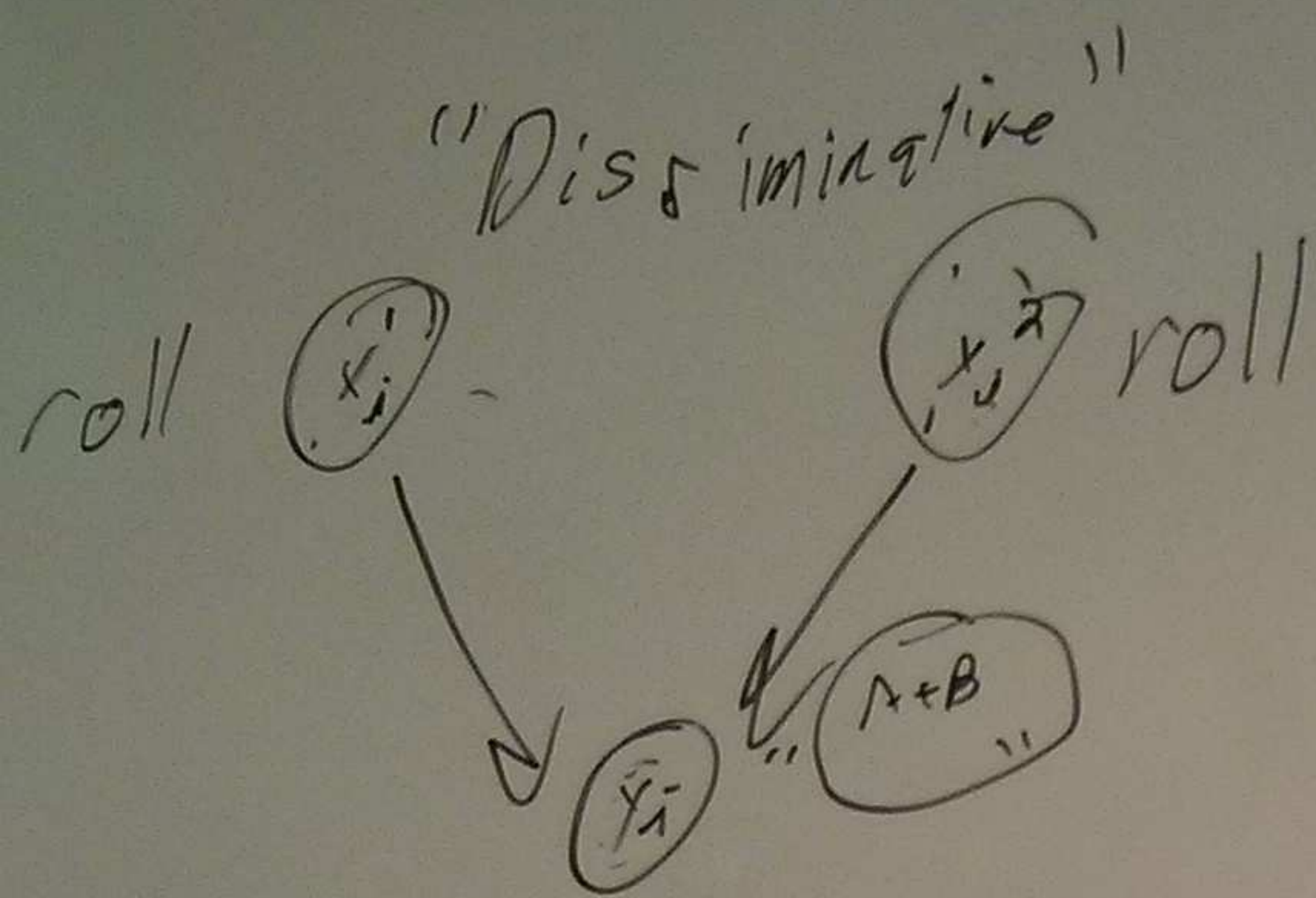
$$p(y_i | x_i) = \frac{p(x_i | y_i) p(y_i)}{p(x_i)}$$

$\propto p(x_i | y_i) p(y_i)$

hard

$x_i^1$	$x_i^2$	$y_i$	
0	0	0	$p(0,0 0)$
0	0	1	$p(0,0 1)$
0	1	0	$\vdots$
0	1	1	$\vdots$
1	0	0	$\vdots$
1	0	1	$\vdots$
1	1	0	$\vdots$
1	1	1	$\vdots$

$2^{2D+1}$  parameters



# Naive Bayes

$x^i$  are "mutually independent" given  $y$

$\forall$  subsets  $A, B$   $x^A \perp x^B \mid y$

$$p(x_{i:1:D} \mid y_i) = p(x_{i:1} \mid x_{i:2:D}, y_i) p(x_{i:2:D} \mid y_i)$$

$$= p(x_{i:1} \mid y_i) p(x_{i:2:D} \mid y_i)$$

$$= \prod_{j=1}^D p(x_{i:j} \mid y_i)$$

↳ "easy"

$x^j$	$y$
0	0
0	1
1	0
1	1

40

ice  $\rightarrow 1, 2, \dots, 6$

# Maximum likelihood

$$\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta) \stackrel{iid}{=} \operatorname{argmax}_{\theta} \prod_{i=1}^N p(D_i | \theta)$$

$$\updownarrow$$

$$\operatorname{argmax}_{\theta} \log p(\mathcal{D} | \theta) \stackrel{iid}{=} \operatorname{argmax}_{\theta} \sum_{i=1}^N \log p(D_i | \theta)$$

numerical

H
T

log

$\theta: p(\text{"head"})$

$$\prod_{i=1}^N p(x_i, y_i) \stackrel{iid}{=} \prod_{i=1}^N p(x_i | y_i) p(y_i)$$

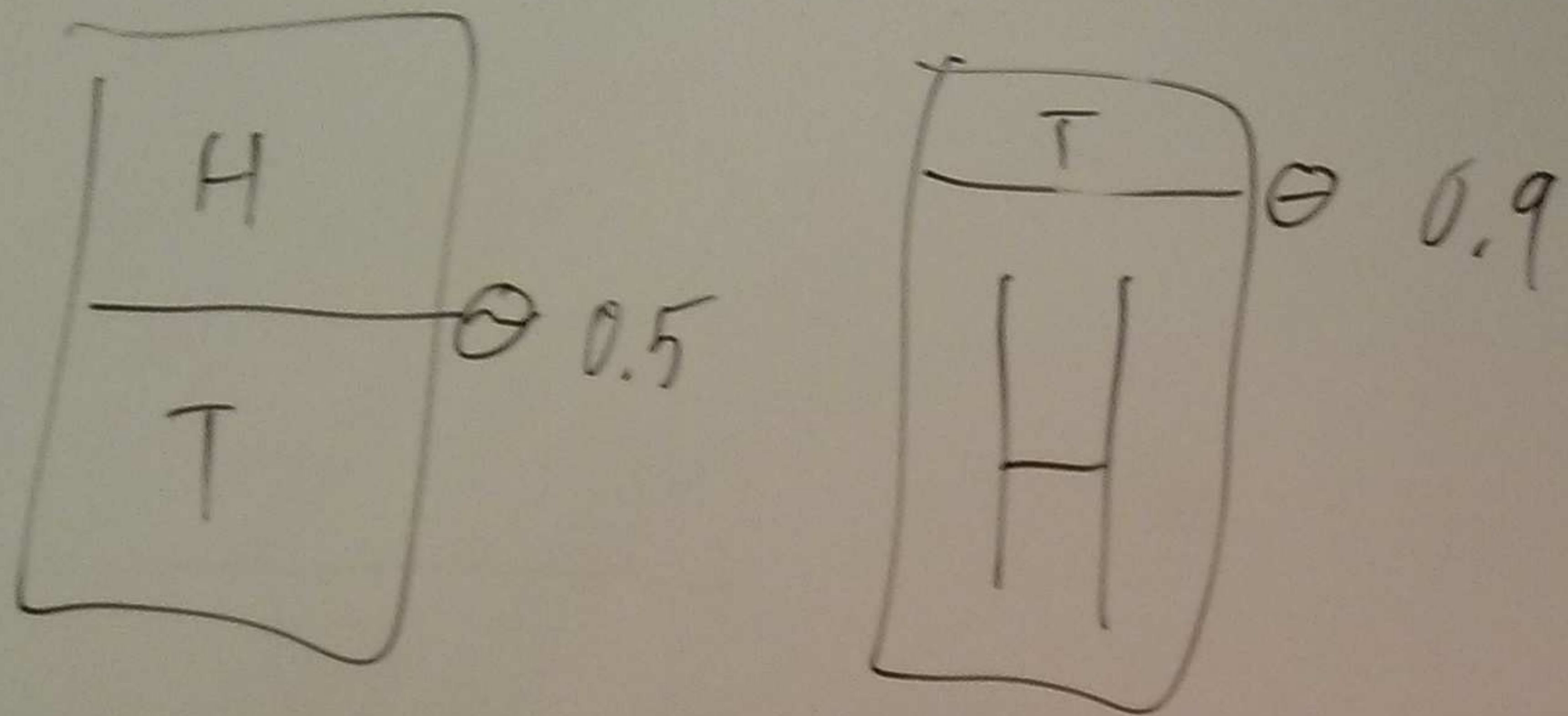
$$= \prod_{i=1}^N p(y_i) \prod_{j=1}^p p(x_i^j | y_i)$$

$$= \prod_{j=1}^p \operatorname{Ber}(y_i | \theta) \prod_{j=1}^p \operatorname{Ber}(x_i^j | y_i, \theta_j^{y_i})$$

$p(y_i)$

	$y$
0	0
1	1
1	0
1	1

40



$p(x_i | \theta)$

naical

$p(x_i | \theta)$

$\theta: p(\text{"head"})$

$$\log p(x, y) = \sum_{i=1}^N \log p(y_i | \theta)$$

$$+ \sum_{j=1}^{\Phi} \log p(x_i^j | y_i, \theta_j^{y_i})$$

$y_i | \theta \sim \text{Ber}(\theta)$  <sup>'heads'</sup>

$p(y_i | \theta) = \theta^{I(y_i=1)} (1-\theta)^{I(y_i=0)}$  <sup>'tails'</sup>  $0 \leq \theta \leq 1$

$$N_1 + N_0 = N$$

$$(1 - \theta)N_1 = N_0 \theta$$

$$\frac{N_1/N}{N_0/N} = \frac{\theta}{1-\theta}$$

$$\theta = \frac{N_1}{N}$$

$$y_i \in \{1, 2, 3, 4\} \quad c=1, 2, \dots$$

$$x_i | y_i, \theta_j^{y_i}$$

$$\theta^c = \frac{N_c}{N}$$

$$\theta_j^c = \frac{N_{jc}}{N_c}$$

$$2(1-\theta)$$