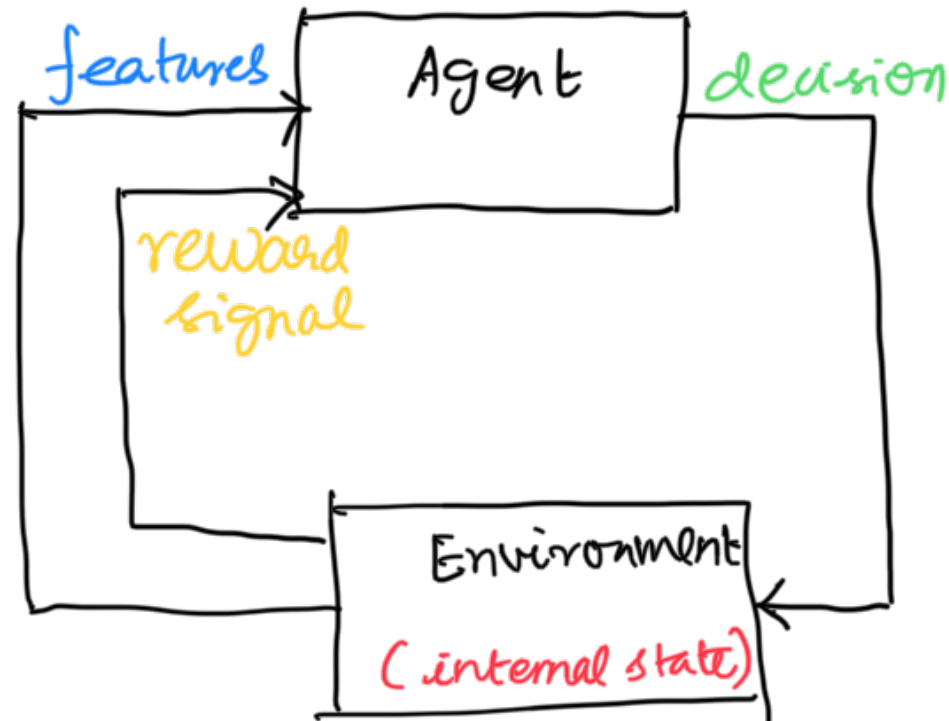


# Statistical Decision Theory (Bird's eye view)



max  
(over decision)

total reward

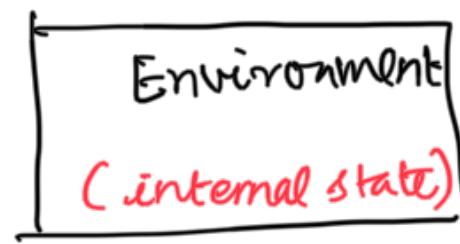
(Decision theoretic  
Goal)

Classification → probabilistic / decision / reinforcement

Statistics → probabilities / frequency of occurrences  
and co-occurrences

## Probability Model (Motivating example)

Say I would like to model an environment/world  
which has 360 days and 20% days it rains.



1 = rain

0 = no rain

$y_t$ : internal state (rain (no-rain))

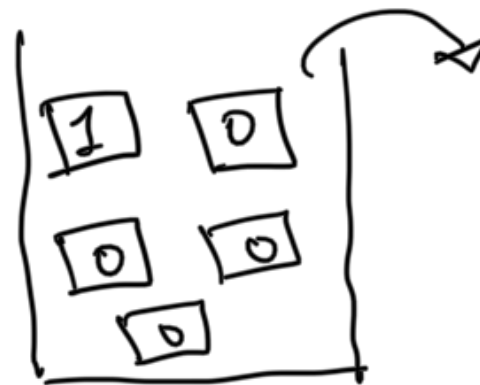
Deterministic 1)  
model

$y_t = 1, 1, \dots, 1, 0, 0, \dots, 0$   
⏟  
72 days
↑  
day 360

Deterministic  
Model 2)

$y_t = 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, \dots$

Stochastic  
(Probabilistic)  
Model

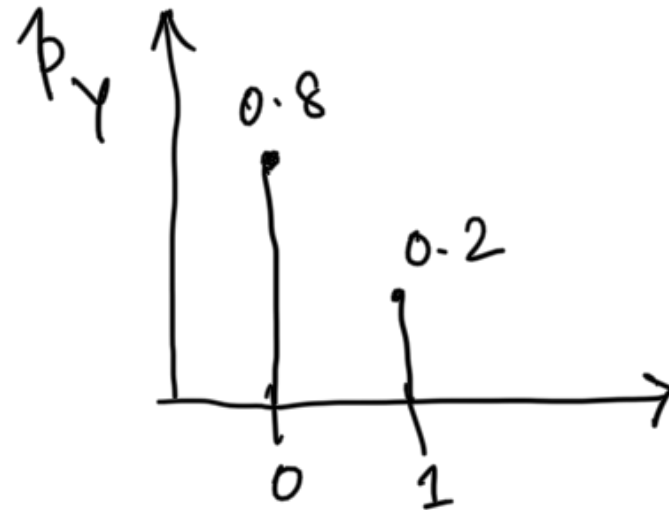


at time  $t$

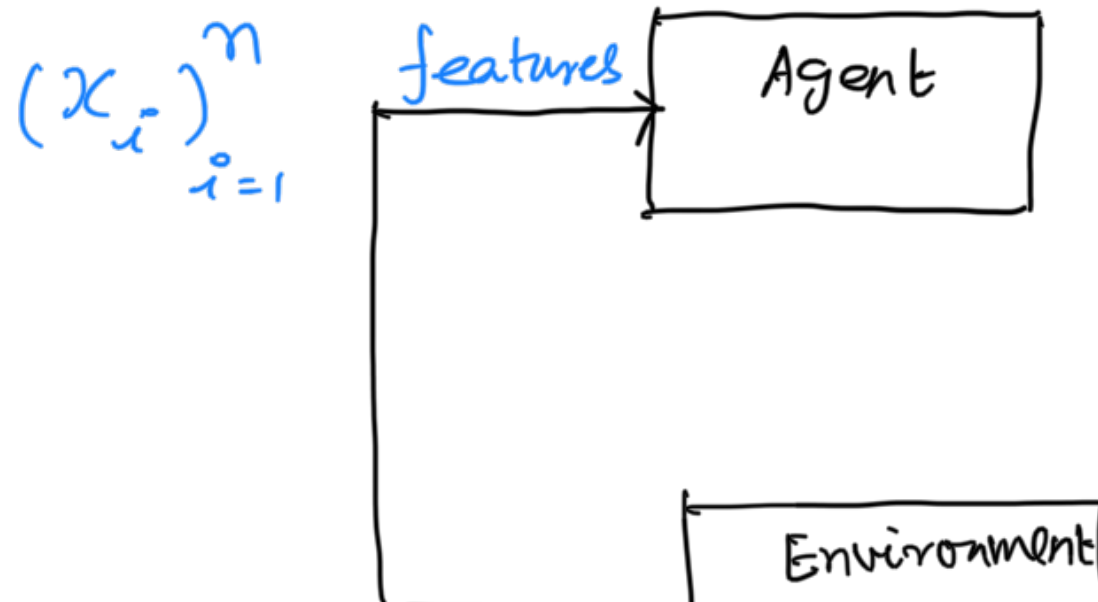
\* shuffle well

\* pick a bit and read out

$y_t \stackrel{i.i.d}{\sim} p_y(\cdot)$ , where



Descriptive Task  
(Data Representation)



(internal state)

Agent does not make any decision

- Features :  $(x_i)_{i=1}^n \in \mathbb{R}^d$

Eg 1) clustering : Goal is to say group articles in a newspaper by topic

$n$ : articles,  $x_i \in \mathbb{R}^d$  ( $d$  = size of vocabulary)

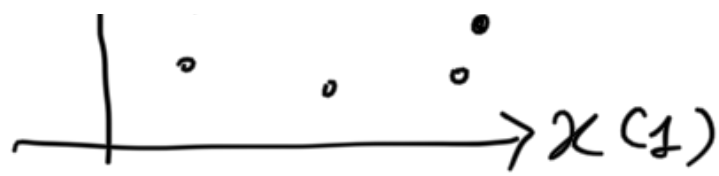
$x_i(j)$  = # word  $j$  occurred in article  $i$ .

$x(2)$



$x(2)$

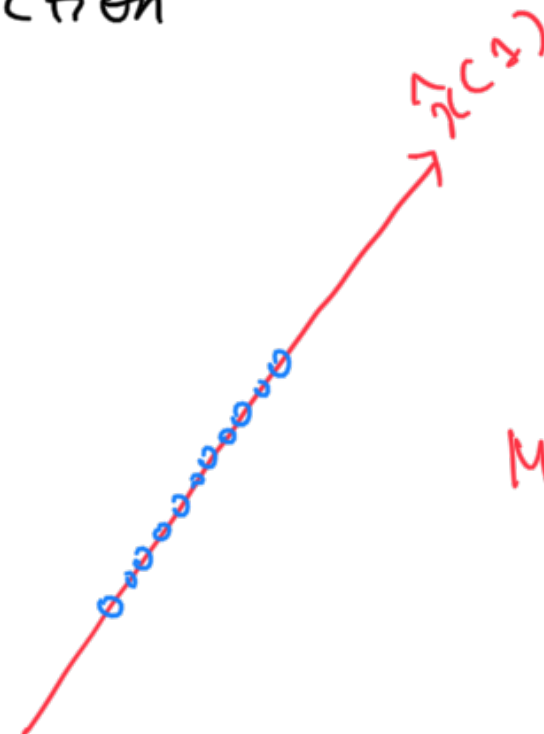
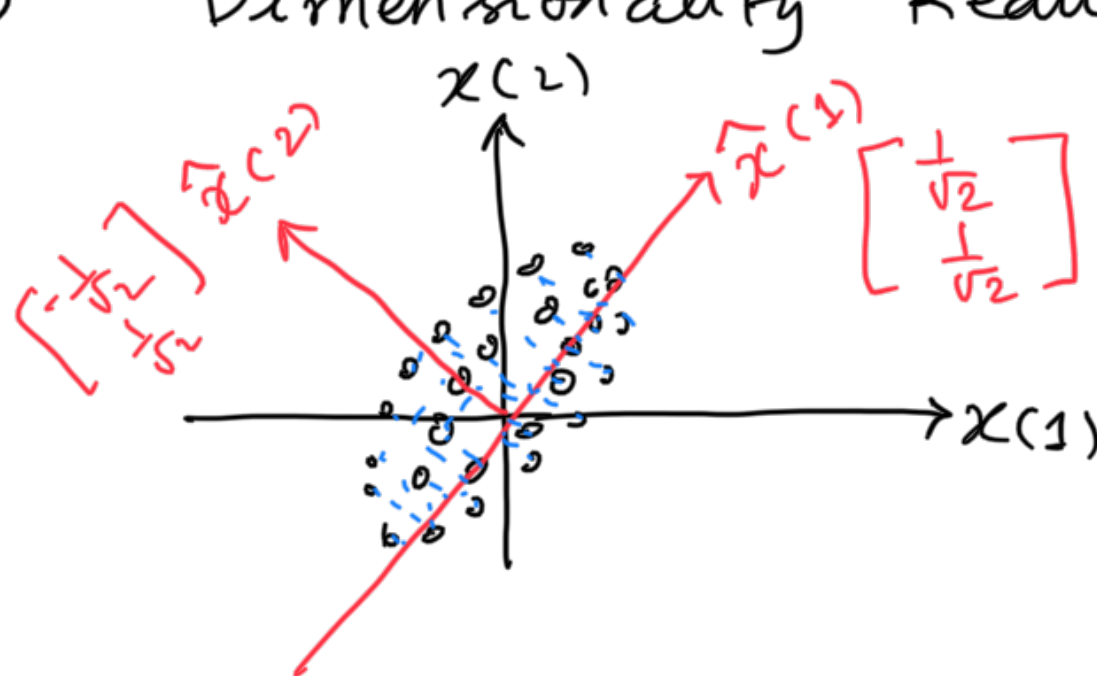




clustering  $\mathcal{C} : \{1, \dots, n\} \rightarrow \{1, \dots, k\}$   
↓  
total clusters

$$x_i \in \text{cluster } \mathcal{C}(i)$$

Eg (2) Dimensionality Reduction



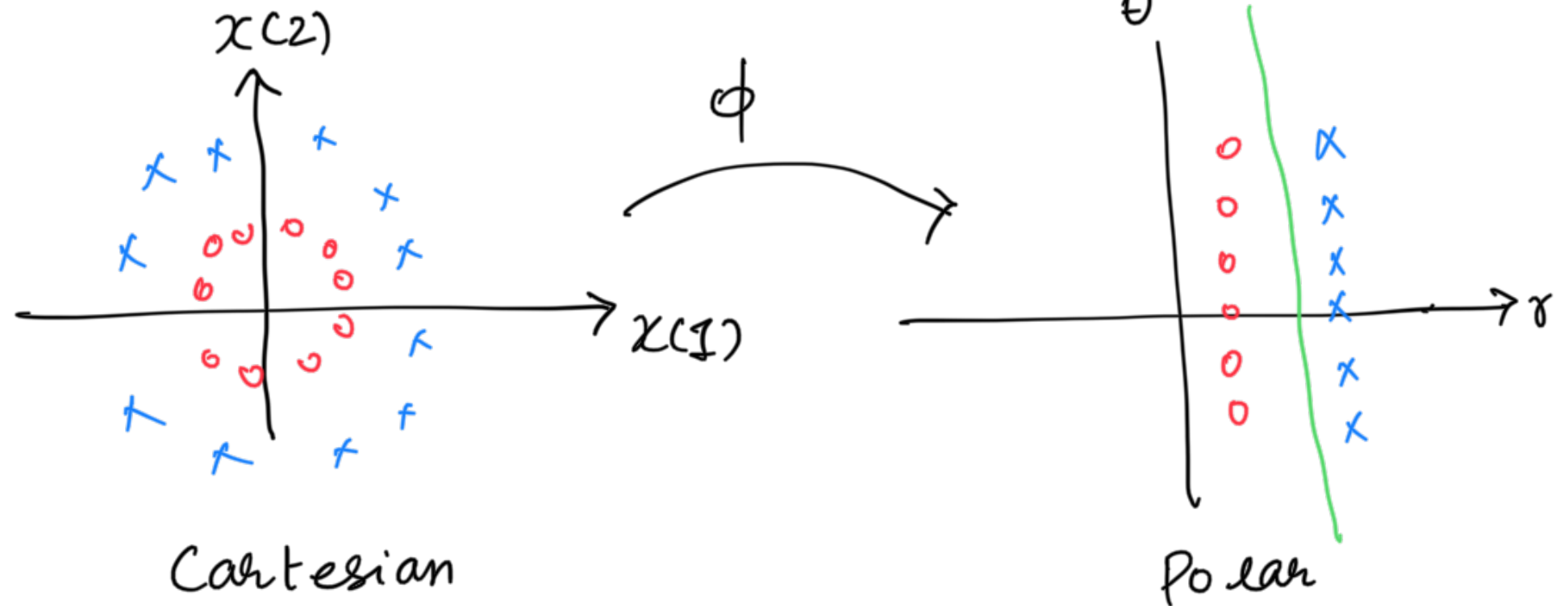
$$M = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \checkmark$$

$$(x_i)_{i=1}^n \in \mathbb{R}^d \xrightarrow{M} (\hat{x}_i)_{i=1}^n \in \mathbb{R}^{d'} \quad (d' \ll d)$$

$$\hat{x}_i = M x_i$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Ex (3) : Representation Learning



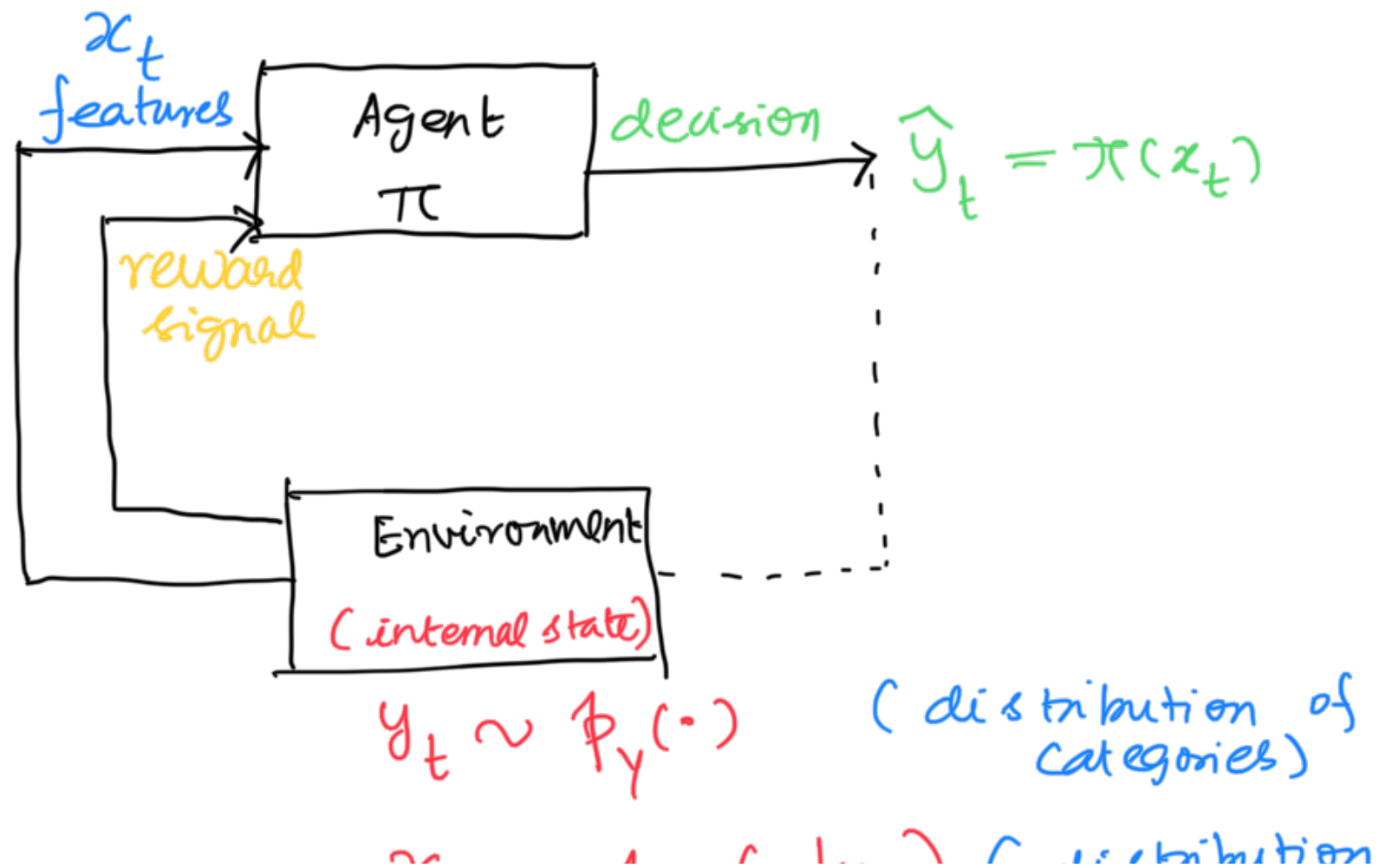
$$x \in \mathbb{R}^d \rightarrow \mathbb{R}^m$$

$$\psi : \mathcal{X} \rightarrow \mathcal{Y}$$

$$\hat{x}_i = \phi(x_i)$$

Predictive Task

Static Prediction





$$x_t \sim p_{x|y}(\cdot | y_t) \quad (\text{distribution of image given category})$$

Multi-class classification (object classification)

- $y_t \in Y = \{1, \dots, c\}$  (label space) → categories

For instance, object classification

$$Y = \{\text{house, person, dog, cat, } \dots, \text{car, plane, } \dots\}$$

- $x_t \in X$  (feature space)  $\subseteq \mathbb{R}^d$

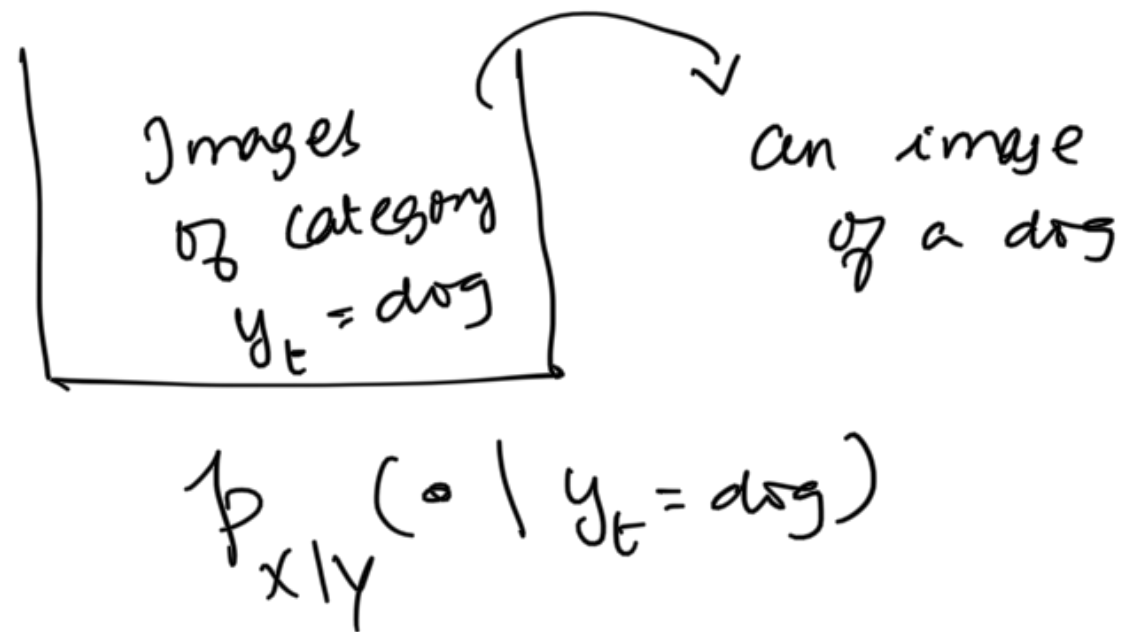
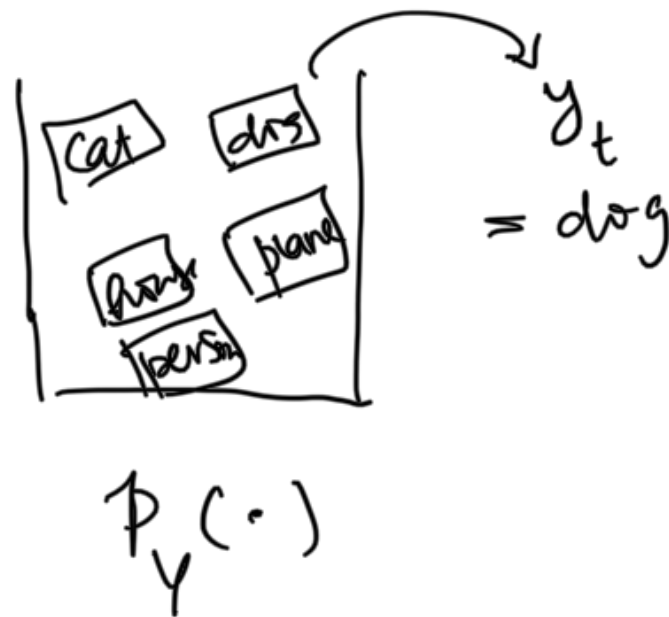
In object classification  $x_t$  is an image

$$x_t \in \mathbb{R}^{\# \text{ pixels}}$$

- $\pi: X \rightarrow Y$

- $l_t = L(\text{state, decision})$   
 $= L(\text{label, predicted label})$

- $L(y_t, \hat{y}_t) = 0, \hat{y}_t = y_t$   
 $= 1, \hat{y}_t \neq y_t$



Goal:  $\min_{\pi} \mathbb{E} \left[ \sum_{t=1}^T l_t \right]$

Model is known, i.e.,  $b, b, \dots$

$$\hat{y}_t = \pi_*(x_t) = \arg \max_{y \in Y} p_Y(y) \cdot p_{X|Y}(x_t | y)$$

(Bayesian Decision Theory)

### Regression

- $X \subseteq \mathbb{R}^d$  (feature)
- $Y \subseteq \mathbb{R}^m$
- $(x_t, y_t) \stackrel{\text{iid}}{\sim} p_{XY}(\cdot, \cdot)$
- $\ell_t = L(y_t, \hat{y}_t)$

$$L(u, \hat{u}) = \|u - \hat{u}\|^2$$

$$\| \partial_t, \partial_t \| = \| \partial_t, \partial_t \|_2$$

$$\text{Goal : } \min_{\pi} \mathbb{E} \left[ \sum_{t=1}^T \ell_t \right]$$

$$\hat{y}_t = \pi_{\star}(x_t) = \mathbb{E}[Y | X = x_t]$$

(Bayesian Decision Theory)

$$(x_t, y_t)_{t=1}^T$$