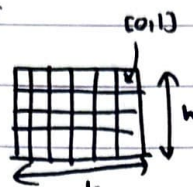
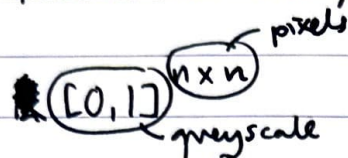


Machine learning reading group, talk 2 (Supervised learning) X = set of inputs Y = set of outputs

① pixellated images



Dog? No, Yes

 $\{0, 1\}$ ② $\{\text{area of house, city, rooms}\} \in \mathbb{R} \times \mathbb{N}^2 \rightarrow \text{Sale price} \in \mathbb{R}$ ③ $\{\text{height, weight, blood pressure}\} \in \mathbb{R}^3 \rightarrow \text{Average steps a day} \in \mathbb{R}$ ~~Deterministic or not? How are X and Y related?~~~~Deterministic or not? How are X and Y related?~~① is deterministic

$\exists f^*: X \rightarrow Y$, "Given an image, the label ~~No or Yes~~ is determined (must be one or other)"

~~No~~ $\hat{0.5}$ dog or 0.6 dog.

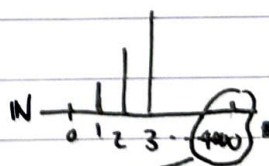
②, ③ are not deterministic

\exists prob dist p on $(X \times Y)$
(think of conditional $y|x$)

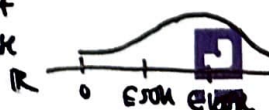
②

 \exists prob dist on X

e.g. have very unlikely to have 4000 rooms



" \exists prob dist on $y|x$ "
= given fixed features x of house, what is prob of house selling at prob y



We just learnt: X, Y, p

Rule: If deterministic, we can get a p from f^* .
Take $y|x$ as $f^*(x)$ with probability 1
(so can always work with p).

However! Data sets in practice don't know p .

"no correlation
and from same
dist"

But, we have n independent samples from p .

independent
and identical
distributed
samples from
 p

Theory: $D_n(p) = \{(x_1, y_1), \dots, (x_n, y_n)\}$
random set that depends on p

Example ②: $D_n(p) =$ set of n random data samples
from land registry data of house
into against sale price

We write $A =$ algorithm, which is a function

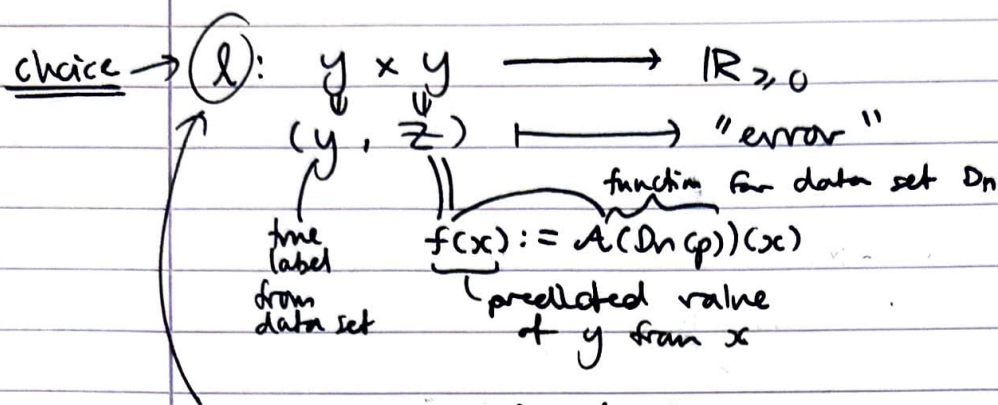
$A : (X \times Y)^{(n)} \longrightarrow \{f: X \rightarrow Y\}$
 \downarrow \uparrow
 $D_n(p)$ # data points
trained on $D_n(p)$

\Rightarrow algorithm returns a deterministic rule
① yes or no decision
exact answer or opposed to
prob of answer
② prob dist of
prices

Example ②: House prices are nondeterministic, but still
useful to have algo ^{that} spit out deterministic
function (good enough for most cases).

②

Fix $p, D_n, \overset{\text{choice}}{\underset{\text{no choice}}{A}}$, we define a loss function



independent of A , "evaluate every algorithm A in same way"

Example ②: take $x = \{\overset{\text{area}}{100 \text{ m}^2}, \overset{\text{city}}{\text{coventry}}, \overset{\text{rooms}}{3}\}$

$\Rightarrow z = f(x) = A(D_n(p))(x) \leftarrow$ predicted value for given x

take $y =$ true value from data set

Define loss function $l(y, z) = \|y - z\|$, e.g.

$l(y, z) = \|y - z\| = 10 \text{K}$ \leftarrow decide if good or not

$\in [0, 100 \text{K}]$ $\in [0, 200 \text{K}]$

Finally we can define risk

$$R(f) := \mathbb{E}[l(y, f(x))] = \text{"average loss over all data"}$$

$f = A(D_n(p))$ $(x, y) \sim p$ dist. over x, y

$$= \int_{x, y} l(y, f(x)) d p(x, y)$$

Risk cannot be computed in practice as we do not have access to p , instead we use empirical risk

$$\hat{R}_n(f) := \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

\uparrow just taking over data $D_n(p)$
(What we know)

$R_{\text{with}}: n \rightarrow \infty$, converge to risk.

Goal of supervised learning:

In practice, fix suitable l , we want to choose an h s.t $\hat{R}_n(f)$ is minimized.

over next four weeks : learn different h 's and in each case how to minimise $\hat{R}_n(f)$
