## Threshold classifier

- hard response, ~~this~~ $X \in \mathbb{R}$ (1-dimensional), $Y \in \{-1, +1\}$

  decision rule
  $$\hat{Y} = \begin{cases} 1 & \text{if } X > x_0 \\ -1 & \text{if } X < x_0 \end{cases} \quad (X = x_0) \text{ whatever})$$
  $$\underset{\uparrow \text{ threshold}}{}$$

- example: $X$ is a person's weight $\quad Y = 1$ obese $\quad Y = -1$ normal

  $x_0 = 100 \, kg$, but this is not a very good method, because it doesn't adjust for a person's height

  ⇒ solution for several features: design a formula that combines the different features into a single number ($\hat{=}$ new feature)

- example: $X \in \mathbb{R}^2$ $\quad X_1 : \underset{[kg]}{\text{weight}}$ $\quad X_2 : \underset{[cm]}{\text{height}}$ $\quad Y = 1$ obese

  body-mass-index $\quad BMI = \dfrac{\text{weight}}{\text{height}^2}$

  decision rule $\quad Y = 1 \quad$ if $\quad BMI \geq 30$

  actually: use multiple thresholds for more informative labels
  $$\hat{Y} = \begin{cases} \text{severe underweight} & \text{if } BMI \leq 16 \\ \text{underweight} & BMI < 18.5 \\ \text{normal} & < 25 \\ \text{over-weight} & < 30 \\ \text{obese} & \geq 30 \end{cases}$$

threshold classifier: given $X \in \mathbb{R}$, predict label $\hat{y} = \begin{cases} 1 & \text{if } X \geq t \\ 0 & \text{if } X < t \end{cases}$

$\Big[$ indicator function: $\mathbb{1}[\text{condition}(x)] = \begin{cases} 1 & \text{if } \text{condition}(x) == \text{true} \\ & \text{either true or false,} \\ & \text{depending on value of } x \quad 0 \quad \text{if } \text{condition}(x) == \text{false} \end{cases}$

in programming: type conversion from boolean to reals
$$\underset{\text{condition}}{\uparrow} \qquad \underset{1 \text{ or } 0}{\uparrow}$$

in C/C++ : double result = (double) (x > t); $\Big]$

rewrite threshold: $\hat{y} = \mathbb{1}[X \geq t]$

compare with Bayesian classifier [$\hat{=}$ best possible)

define a toy problem : - simplified "cartoon" of some real world problem
- is not necessarily practically relevant
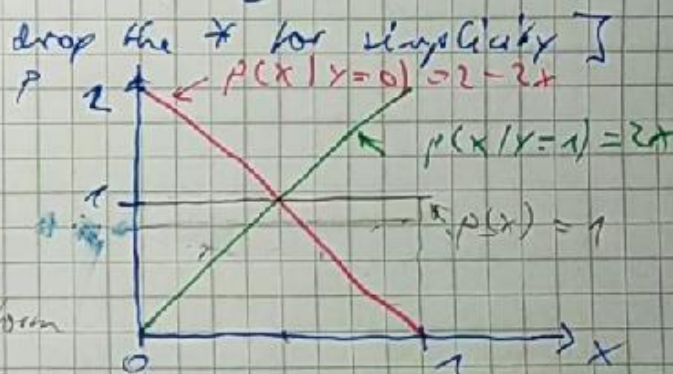- but: we learn a lot about the machine learning
method that we use

here: $X \in [0, 1]$   prior of $Y$ :   $p(Y=1) = p(Y=0) = \frac{1}{2}$
[ this notation actually is $p^*(Y)$, but we drop the $*$ for simplicity ]

likelihoods                    evidence :

$p(X | Y = 0) = 2 - 2x$ $\Big|$ $p(x) = \sum_{u=0,1} p(x|Y=u) p(Y=u)$

$p(x | Y = 1) = 2x$ $\qquad = (2-2x) \cdot \frac{1}{2} + 2x \cdot \frac{1}{2}$

$\qquad\qquad\qquad = 1 - x + x = \underline{\underline{1}}$ uniform

posteriors according to Bayes:  $\quad p(Y|x) = \dfrac{p(X|Y)\,p(Y)}{p(X)}$

$$p(Y=0\,|\,x) = \frac{(2 - 2x)\cdot \frac{1}{2}}{1} = 1-x \qquad p(Y=1\,|\,x) = \frac{2x\cdot \frac{1}{2}}{1} = x$$

Bayes classifier:  $\qquad y^* = \arg\max_{u} p(Y=u\,|\,x)$

$$= \begin{cases} 0 & \text{if } (1-x) \geq x \iff x \leq \frac{1}{2} \\[2mm] 1 & \text{if } (1-x) < x \iff x > \frac{1}{2} \end{cases}$$

$$= \mathbb{1}\left[\,x > \tfrac{1}{2}\,\right] \,\hat{=}\, \text{threshold classifier with } t = \tfrac{1}{2}$$

compare with all possible threshold classifiers:

- type A: $\qquad \hat{Y} = \mathbb{1}[\,x > t\,]$ $\qquad$ - type B: $\quad \hat{Y} = \mathbb{1}[\,x < t\,]$

compute the probability of error $\qquad p(\text{error}\,|\,\text{type},\,t)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ two degrees of freedom

$$p(\text{error}\,|\,\text{type},t) = \mathbb{E}_{\substack{x\sim p(x|Y) \\ Y\sim p(y)}}\left[\,p\left(\,\mathbb{1}[\text{condition}(x)] \neq Y\;|\;\text{type},\,t\right)\,\right]$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\mathbb{1}[\,x > t\,]$ if type=A
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\mathbb{1}[\,x < t\,]$ if type=B

$\left.\begin{array}{l}\text{since } p(Y=1) = p(Y=0) = \frac{1}{2} \\ \text{for classifier A}\end{array}\right\} = \mathbb{E}_x\left[\,p(Y=1|x)\,\mathbb{1}[x < t]\,\right] + \mathbb{E}_x\left[\,p(Y=0|x)\,\mathbb{1}[x > t]\,\right]$

$\qquad\qquad\qquad\qquad\qquad\qquad$ error, since
$\qquad\qquad\qquad\qquad\qquad\qquad$ opposite of type A
$\qquad\qquad\qquad\qquad\qquad\qquad$ behavior

$$E_X\left[p(Y=1|X)\,\mathbb{1}[x<t]\right] = \int_0^1 p(Y=1|x)\,\underbrace{\mathbb{1}[x<t]}_{\text{from } \mathbb{E}}\,\underbrace{p(x)}_{}\,dx$$

false negative

$$\underbrace{\qquad\qquad}_{\text{argument of } \mathbb{E}}$$

$$= \int_0^t \underbrace{p(Y=1|x)}_{=x}\,\underbrace{p(x)}_{=1}\,dx = \int_0^t x\cdot 1\,dx = \frac{x^2}{2}\Big|_0^t = \frac{t^2}{2}$$

$$E_X\left[p(Y=0|X)\,\mathbb{1}[x\geq t]\right] = \int_0^1 p(Y=0|x)\,\mathbb{1}[x\geq t]\,p(x)\,dx$$

false positive

$$= \int_t^1 \underbrace{p(Y=0|x)}_{=1-x}\,\underbrace{p(x)}_{=1}\,dx = \int_t^1 (1-x)\,dx = \left(x-\frac{x^2}{2}\right)\Big|_t^1$$

$$= 1-\frac{1^2}{2} - \left(t-\frac{t^2}{2}\right) = \frac{1}{2} - t + \frac{t^2}{2} =$$

$$E_X\left[p(\text{error}\,|\,\text{type}=A,\,t)\right] = \frac{t^2}{2} + \frac{1}{2} - t + \frac{t^2}{2} = t^2 - t + \frac{1}{2} = \left(t-\frac{1}{2}\right)^2 + \frac{1}{4}$$



$p(Y=0|x)$    $p(Y=1|x)$

type B     type A

Minimum is achieved at $t=\frac{1}{2}$, as in Bayes classifier

Same calculation for **type B** (reverse all conditions)

$$E_X\left[p(\text{error}\,|\,\text{type}=B,\,t)\right] = -\left(t-\frac{1}{2}\right)^2 + \frac{3}{4}$$

$$= 1 - E_X\left[\text{type A}\right]$$

Bayes error: set threshold where the curves cross

Minimum is achieved for $t=0$ or $t=1$

$$E_X\left[p(\text{error})\right] = \frac{1}{2} \;\hat{=}\; \text{pure guessing}$$

in contrast, the best error for type A at $t=\frac{1}{2}$ is $\frac{1}{4}$ [$\hat{=}$ Bayes rate]

likelihoods are probability <u>densities</u>, because $X$ is continuous

normalisation $\int_0^1 p(x|y)\, dx = 1$ for all $y$ [ here: $y=0$ or $y=1$ ]

$$\int_0^1 p(x|Y=0)\, dx = \int_0^1 (2-2x)\, dx = 2x - \frac{2x^2}{2}\Big|_0^1 \overset{?}{=} 2 - 1 - (0-0) = 1$$

$$\int_0^1 p(x|Y=1)\, dx = \int_0^1 2x\, dx = \frac{2x^2}{2}\Big|_0^1 = 1 - 0 = 1$$

---

How to generalize the threshold classifier when there are multiple features

$X \in \mathbb{R}^D$ $\Rightarrow$ should be better, because more features $\hat{=}$ more information on $Y$

problem : comparison $X \geq t$ is only defined for scalars $X \in \mathbb{R}$, not vectors

solutions : ① reduce $X$ to a 1-dimensional score : $z = g(X) \in \mathbb{R} \Rightarrow \hat{Y} = \mathbb{1}[z \geq t]$

~~problem~~ example : body-mass-index : $t = BMI = \dfrac{\text{weight}}{\text{height}^2} = \dfrac{X_0}{X_1^2}$

$[x_0 : \text{~~weight~~ weight}, x_1 : \text{height}]$

problem : finding good function $g(x)$ is <u>hard</u> $\Rightarrow$ learn $g(x)$ $\Rightarrow$ later

② reduce multi-dimensional comparison to a sequence of 1-dimensional comparisons over the elements of $X$ ($\hat{=}$ single features) $\Rightarrow$ decision tree $\Rightarrow$ later

③ nearest neighbor classifier : define the threshold implicitly via distances to "representatives" for every class $\Rightarrow$ next