

$$p(\underbrace{\text{error}}_{\hat{y} = \hat{f}(x) \neq y^*} \mid \text{type} = A, t) \left. \begin{array}{l} \leftarrow \text{some fixed threshold} \\ \text{and } \hat{y} = \mathbb{1}[x \geq t] \end{array} \right\} \text{threshold classifier error}$$

expected error: average error over all possible test data points

$$\mathbb{E}_{x, y^* \sim p^*(x, y)} [p(\hat{f}(x) \neq y^* \mid \text{type} = A, t)]$$

$$\sum_{k=0}^1 \int p^*(x, y=k) p(\hat{f}(x) \neq k \mid \text{type} = A, t) dx = \left(t - \frac{1}{2}\right)^2 + \frac{1}{4}$$

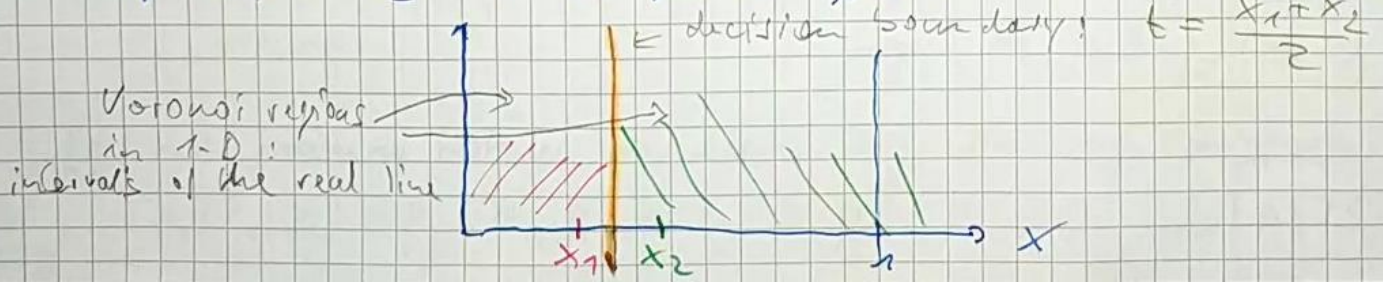
$\Rightarrow \mathbb{E} \Rightarrow \sum$ \uparrow classes are discrete
 $\mathbb{E} \Rightarrow \int$ \uparrow features are continuous

apply this result to the nearest-neighbor classifier for $N=2$

- we do not fix t manually, but by means of the TS: $t = \frac{x_1 + x_2}{2}$

$$\{(x_1, y_1=0), (x_2, y_2=1)\}$$

if $x_1 < x_2 \Rightarrow$ we get a type A threshold classifier, $x_1 > x_2 \Rightarrow$ type B



to compute the expected error of the nearest neighbor classifier for $N=2$, we must take two averages:

- average about all possible test points (as before)
- average over all possible training sets, because the classifier performance depends on the training set

e.g. $TS_1 = \{ (x_1=0.2, y_1=0), (x_2=0.8, y_2=1) \}$
 $t = \frac{0.2+0.8}{2} = 0.5, x_1 < x_2 \Rightarrow \text{type A}$
 good

$TS_2 = \{ (x_1=0.5, y_1=0), (x_2=0.4, y_2=1) \}$
 $t = \frac{0.5+0.4}{2} = 0.45, x_1 > x_2 \Rightarrow \text{type B}$
 bad

$$E_{TS} \left[\frac{p(\text{error} | TS)}{1} \right] = p(\text{error} | \text{type}, t) \text{ with type and } t \text{ calculated from } TS$$

in general $\int p(TS) \cdot p(\text{error} | TS) dTS$

i.i.d. assumption: independent, identically distributed

TS has two properties: ① each training instance was obtained independently of the other training instances

dice: i -th throw outcome independent of $i' < i$ outcomes

$$\Rightarrow p(TS) = p(\{ (x_i, y_i) \}_{i=1}^N) \text{ factorizes} \\ = p_1(x_1, y_1) \cdot p_2(x_2, y_2) \cdot \dots \cdot p_N(x_N, y_N)$$

② each training instance was drawn from the same (unknown) true probability distr.

for all i : use the same dice
 $p_1(x_1, y_1) = p_2(x_2, y_2) = \dots = p^*(x, y) \Rightarrow p(TS) = \prod_{i=1}^N p^*(x_i, y_i)$

$$\begin{aligned} \mathbb{E}_{TS} [p(\text{error} | TS)] &= \int p(TS) p(\text{error} | TS) dTS \\ &\stackrel{\text{i.i.d.}}{=} \int p(x_1 | y_1) \cdot p(x_2 | y_2) \cdots p(x_n | y_n) p(\text{error} | \{(x_i, y_i)\}_{i=1}^n) d(x_1, y_1) \cdots d(x_n, y_n) \end{aligned}$$

specifically for the toy example: we already know that $y_1 = 0, y_2 = 1$

\Rightarrow the integration and prob of y can be dropped

$$p(TS) = p(x_1 | y_1 = 0) \cdot p(x_2 | y_2 = 1) \quad \text{i.i.d. assumption applied to } TS \{ (x_1, y_1 = 0), (x_2, y_2 = 1) \}$$

if in addition, $x_1 < x_2 \Rightarrow$ we get a type A classifier

$$\begin{aligned} \mathbb{E}_{TS} [p(\text{error} | TS)] &= \int_0^1 \int_0^1 p(x_1 | y_1 = 0) \cdot p(x_2 | y_2 = 1) \mathbb{1}[x_1 < x_2] p(\text{error} | \text{type A}, t = \frac{x_1 + x_2}{2}) dx_1 dx_2 \\ &= \int_0^1 \underbrace{p(x_1 | y_1 = 0)}_{2-2x} \left[\int_{x_1}^1 \underbrace{p(x_2 | y_2 = 1)}_{2x} \underbrace{p(\text{error} | \text{type A}, t = \frac{x_1 + x_2}{2})}_{\left(t - \frac{1}{2}\right)^2 + \frac{7}{4} = \left(\frac{x_1 + x_2}{2} - \frac{1}{2}\right)^2 + \frac{7}{4}} dx_2 \right] dx_1 \end{aligned}$$

\Rightarrow give all this to wolfram (cloud) and get $\frac{83}{360}$

$$\mathbb{E}_{TS} [p(\text{error} | TS)]_{x_2 < x_1} = \frac{43}{360} \quad p(TS | y_1 \neq y_2)$$

$$\mathbb{E}_{TS} [p(\text{error} | TS)] = \mathbb{E}_{TS}^{x_1 < x_2} + \mathbb{E}_{TS}^{x_1 > x_2} = 35\%$$