

# Classifier Design Based on Expected Risk (Loss/Cost) Minimization

Assume that our data samples  $\{x_1, x_2, \dots\}$  are independent identically distributed (iid) instances of a real-valued  $n$ -dimensional random vector  $X \in \mathbb{R}^n$ .

Let the probability <sup>density</sup> function (pdf) of  $X$  be a mixture (convex linear combination) of  $C$  class-conditional pdfs as follows:

$$P_X(x) = p_{X|L}(x|1)p_L(1) + \dots + p_{X|L}(x|C)p_L(C)$$

Notation of Papoulis  $\rightarrow$

$$= p(x|L=1)p(L=1) + \dots + p(x|L=C)p(L=C)$$

Less precise but more convenient notation  $\rightarrow$

$$= \sum_{l=1}^C \underbrace{p(x|L=l)}_{\text{Data pdf for class label } l} \underbrace{p(L=l)}_{\text{Prior probability for class label } l}$$

CLASS-CONDITIONAL PRIOR  $\rightarrow$

Suppose that we want to design a decision rule (classifier) that selects the "best" choice among a discrete set of options for a given  $x \in \mathbb{R}^n$ :  $D: \mathbb{R}^n \rightarrow \{1, \dots, A\} \subset \mathbb{Z}$

In expected risk (loss/cost) minimization, the "best" option

for convenience  
assume options are  
indexed by integers.

is to select the one which results in the smallest expected risk.

Let  $\Lambda = \begin{bmatrix} \lambda_{11} & \dots & \lambda_{1c} \\ \vdots & & \vdots \\ \lambda_{A1} & \dots & \lambda_{Ac} \end{bmatrix}$  be the loss matrix, such

that  $\lambda_{ij}$  = "the risk/loss/cost associated with deciding on option  $i$ , given  $x$  comes from class label  $j$ "

Note that  $i \in \{1, \dots, A\}$  and  $j \in \{1, \dots, C\}$ .

That is, in general, the set of decision options do not need to be the same as the set of class labels.

We choose  $\lambda_{ij} \geq 0$  for all  $(i, j)$  pairs.

We want to find a decision rule  $D$  that minimizes expected risk/loss/cost. Overall expected risk is as follows:

$$E_X[\text{Risk}] = \int_{-\infty}^{\infty} \underset{\text{Risk}(D(x)=\cdot|x)}{\text{Risk}(D(x)|x)} P_X(x) dx$$

Note that we have only options  $D(x) \in \{1, \dots, A\}$

for a given  $x$ . The risk of deciding  $D(x)=d$  for a given  $x$  is as follows:

$$\text{Risk}(D(x)=d|x) = \sum_{l=1}^C \lambda_{dl} \overbrace{P(L=l|x)}^{\text{class posterior for label } l}$$

the summation here produces the average loss of deciding  $d$  given  $x$ , considering the class posterior probabilities for the given  $x$ .

loss we would incur by deciding  $d$  given sample from class  $l$  !  
 Probability of class label being  $l$  given  $x$

Since  $\lambda_{dl} \geq 0$  and  $P(L=l|x) \geq 0$  for all  $(d,l)$  pairs,

$\text{Risk}(D(x)=d(x)) \geq 0$  for all  $d \in \{1, \dots, A\}$ .

Back to  $E_{\mathbf{X}} \{ \text{Risk} \} = \int_{-\infty}^{\infty} \text{Risk}(D(x)=\cdot | x) P_{\mathbf{X}}(x) dx$

we notice now that both  $\text{Risk}(D(x)=\cdot | x) \geq 0$  and  $P_{\mathbf{X}}(x) \geq 0$  for all  $x \in \mathbb{R}^n$ . Also, we notice that  $\text{Risk}(D(x)=\cdot | x)$  only depends on  $x$  (and not on decisions we make for other  $x'$ ).

Consequently, if we make the minimum-risk decision for each individual  $x$ , we will end up minimizing  $E_{\mathbf{X}} \{ \text{Risk} \}$  overall.

So, the expected-risk-minimization (ERM) decision rule becomes:

$$D(x) = \arg \min_{d \in \{1, \dots, A\}} \text{Risk}(D(x)=d | x)$$

$$= \arg \min_{d \in \{1, \dots, A\}} \sum_{l=1}^C \lambda_{dl} P(L=l | x)$$

$$= \arg \min_{d \in \{1, \dots, A\}} \sum_{l=1}^C \lambda_{dl} \frac{P(x | L=l) P(L=l)}{P(x)}$$

Bayes Rule

$$= \arg \min_{d \in \{1, \dots, A\}} \sum_{l=1}^C \lambda_{dl} P(x | L=l) P(L=l)$$

multiply by  $P(x)$ , which does not change the decision

Basically, for a given  $x$ , we will compute

$$R(D=1|x) = \sum_{l=1}^C \lambda_{1l} P(L=l|x)$$

$$R(D=2|x) = \sum_{l=1}^C \lambda_{2l} P(L=l|x)$$

$\vdots$

$$R(D=A|x) = \sum_{l=1}^C \lambda_{Al} P(L=l|x)$$

class posteriors

and choose the decision with the smallest  $R(D(x)=\cdot|x)$  value as the best option.

Note that, in vectorized form, these equations can be written as

$$\begin{bmatrix} R(D=1|x) \\ \vdots \\ R(D=A|x) \end{bmatrix} = \underset{\substack{\uparrow \\ \text{loss} \\ \text{matrix}}}{\Lambda} \begin{bmatrix} P(L=1|x) \\ \vdots \\ P(L=C|x) \end{bmatrix}$$

$\nwarrow$  all decision risks across options
 $\nwarrow$  all class posteriors

As shown before, the class posteriors can be computed from class conditionals and class priors:

$$P(L=l|x) = \frac{P(x|L=l)P(L=l)}{p(x)}$$

where  $p(x) = \sum_{j=1}^c P(x|L=j)P(L=j)$

The class-posterior vector is, then

$$\begin{bmatrix} P(L=1|x) \\ \vdots \\ P(L=c|x) \end{bmatrix} = \frac{1}{p(x)} \begin{bmatrix} P(x|L=1)P(L=1) \\ \vdots \\ P(x|L=c)P(L=c) \end{bmatrix}$$

$$= \frac{1}{p(x)} \begin{bmatrix} P(L=1) & & 0 \\ & \ddots & \\ 0 & & P(L=c) \end{bmatrix} \begin{bmatrix} P(x|L=1) \\ \vdots \\ P(x|L=c) \end{bmatrix}$$

Later we will see discriminative models that attempt to approximate class posteriors directly, and generative models that attempt to approximate class conditionals/priors.

## Special Case: Minimum Probability of Error Classification Rule

Let  $l, d \in \{1, \dots, C\} \Leftrightarrow$  our decisions will come from the same set as class labels

$$\text{Let } \lambda_{dl} = \begin{cases} 0 & \text{if } d=l \\ 1 & \text{if } d \neq l \end{cases} \quad \begin{array}{l} 0 - \text{loss for correct decision} \\ 1 - \text{loss for any incorrect decision} \end{array}$$

$(\delta_{dl} = \text{Kronecker delta}) \Rightarrow \lambda_{dl} = 1 - \delta_{dl}$

This is called the 0-1 loss and it indicates our design choice that all incorrect decisions are equally bad and they are worse than all correct decisions which are equally good.

$$R(D=d|x) = \sum_{l=1}^C \lambda_{dl} P(L=l|x)$$

$$\lambda_{dl} = 1 - \delta_{dl} \downarrow = \sum_{\substack{l=1 \\ l \neq d}}^C P(L=l|x)$$

$$= 1 - P(L=d|x)$$

After this substitution, the ERM decision rule simplifies to the Maximum A Posteriori (MAP) decision rule:

$$D(x) = \operatorname{argmin}_{d \in \{1, \dots, C\}} R(D = d | x)$$

$$= \operatorname{argmin}_{d \in \{1, \dots, C\}} 1 - P(L = d | x)$$

$$= \operatorname{argmax}_{d \in \{1, \dots, C\}} P(L = d | x)$$

For a given  $x$ , decide on the class label with the largest posterior probability.

MAP classification rule (i.e. ERM with 0-1 loss) achieves minimum probability of error overall.



## Special Case of MAP Classifier: ML Classifier

If the class priors are equal, then MAP classifier simplifies to ML classifier  
ML: maximum likelihood.

$$D_{\text{MAP}}(x) = \arg \max_{d \in \{1, \dots, C\}} P(L=d|x)$$

$$= \arg \max_d \frac{P(x|L=d) P(L=d)}{p(x)}$$

Multiply  
with  $p(x)$

$$= \arg \max_d P(x|L=d) P(L=d)$$

$P(L=d) = \frac{1}{C}$   
for all  $d$   
multiply  
with  $C$

$$= \arg \max_{d \in \{1, \dots, C\}} P(x|L=d)$$

select the class label  
that has a class cond-  
pdf which makes  $x$   
most likely.

If you want to minimize probability of error, use 0-1 loss ( $\Rightarrow$  MAP). In that case if class priors are equal, you can simply use this ML classification rule.

Special Case: ERM with 2 classes.

Let  $d \in \{0, 1\}$ . Then

$$D_{\text{ERM}}(x) = \underset{d \in \{0, 1\}}{\operatorname{argmin}} R(D=d|x)$$

$$= \underset{d \in \{0, 1\}}{\operatorname{argmin}} \lambda_{d0} P(L=0|x) + \lambda_{d1} P(L=1|x)$$

i.e. we will decide as follows

$$\lambda_{00} P(L=0|x) + \lambda_{01} P(L=1|x) \underset{D(x)=0}{>} \underset{D(x)=1}{<} \lambda_{10} P(L=0|x) + \lambda_{11} P(L=1|x)$$

↓ let  $\lambda_{01} - \lambda_{11} > 0$  and  $\lambda_{10} - \lambda_{00} > 0$   
(i.e. incorrect decisions are worse than correct decisions)

Then

$$\frac{P(L=1|x)}{P(L=0|x)} \underset{D(x)=0}{>} \underset{D(x)=1}{<} \frac{(\lambda_{10} - \lambda_{00})}{(\lambda_{01} - \lambda_{11})} \quad \begin{array}{l} \text{Decide based on} \\ \text{the ratio of} \\ \text{class posteriors} \end{array}$$

↓ with Bayes rule |||  $\hookrightarrow$  (if 0-1 loss threshold = 1)

$$\frac{P(x|L=1)}{P(x|L=0)} \underset{D(x)=0}{>} \underset{D(x)=1}{<} \frac{(\lambda_{10} - \lambda_{00}) P(L=0|x)}{(\lambda_{01} - \lambda_{11}) P(L=1|x)} \quad \begin{array}{l} \text{Decide based on} \\ \text{the ratio of} \\ \text{likelihoods} \\ \text{under class cond pdfs.} \end{array}$$