



**Calibration** → model that provides reliable confidence estimates, that allow us to make optimal decision in real world applications

↳ Our model gives us score!

↳ class posteriors → probability that tell us how likely it is that a given point belong to each class

How we know if a model is well calibrated?

act DCF min DCF  
cost using our decision threshold cost with optimal threshold  
well calibrated if these values are close

**Discriminative classifier**

↳ learn distinguish class  
↳ no learn data distribution  
↳ output posterior probabilities

**Generative Classifier**

↳ model data distribution

**objective**

assign higher score to sample from correct class & lower to sample from incorrect class

choose appropriate thresholds for making decision

calibration is here

Poor calibration

↳ good discriminant ability but poorly calibrated scores  
↳ sensitive to prior probability  
no need to retrain model

Model without clear probabilistic interpretation

↳ SVM → score that represent distance

**Over-regularized Underconfident**

↳ much regularization to combat overfitting

QUESTO FA SCUPO A MATEMATICA

CALCOLA ENTRO UN MARGINE CHE PUÒ COPERCHI OTTIME PUNZONI MA NON APPREZZARE → PUÒ SGETTARSI DI NEGLI

Overfit → Overconfident model → overfit to train data → unrealistic posterior probabilities

↳ output value alone to 0 or 1

↳ even in difficult case

SE TI UN VALORE ANTEVVISO CHE PUNTA ALLA CLASSE A, MA NON CONFERMA IL TUE CONFESSIONI, ALLORA IL COTTO EPO (UN PO' SE LO CONFERMI MA NON SE LO CONFIRMI) HA UN PESO PIÙ ALTO. SE TI UN VALORE ANTEVVISO CHE PUNTA ALLA CLASSE B, MA NON CONFERMA IL TUE CONFESSIONI, ALLORA IL COTTO EPO (UN PO' SE LO CONFERMI MA NON SE LO CONFIRMI) HA UN PESO PIÙ ALTO.

↳ push probabilities toward 0.5 even when we have strong evidence for one class

**Distribution Mismatch**

↳ data used for training differs systematically from the reality

↳ Misclassification  
Poor discrimination

For example the gaussian model for the project did is able to well discriminate the sample of two classes orange in this way.

As such it's having a hard time doing the discrimination because around HIGH DENSITY REGION CONF ID 0.55 WE GET A LOT OF DATA, ONLY FEW OF DATA IN THE HIGH DENSITY REGION ARE ORANGE

QUESTO AVVIENE IN UN GENERATIVE MODEL CHE NON È PIÙ POLARE A THIS CLASSIFICATION

How we select a good threshold → We can treat the decision threshold as a hyperparameter that can be optimized in validation set

PROBLEMI DIVERSI CON IL NOSTRO RR (recognizer)

$X_t \rightarrow S(t) \leq t \rightarrow$  CUE INSERITO USATO.

SE ONUFOLIO THRESHOLD  
FA CALARE IL PERC. I.  
NOMI DATI NON SONO  
CALIGRATI BENE, NE  
TROVANDO UNO CHE  
LAVORA MELLO

USATO minDCF  
TRA TUTTI NOI IL  
TASSO DI CUE  
UN PARENTHO  
QUA, UNA VERA

QUOTI PROBLEMI  
[Training set] [Val set] → [App Data]

minDCF → (1) IL TASSO DI CUE  
DEVE ESSERE > QUESTO

IL TASSO DI CUE  
IN TERMINE DI PERFORMANCE  
XII VALUTAZIONE DEI  
L'VALORE DI APP

→ SE IL VALIDATION SET

→ 3 MILLE X L'APPALITZIONE  
SET ALTRUI VA BENE

L'VALORE DI APP

Fragility Problem

poorly calibrated  
score are sensitive

mis-calibrated score

↳ a small shift in  
the prior probability might  
require a completely different  
threshold

Solution? We change approach

We want well-calibrated score from the start so that

Predictable Thresholds

optimal threshold →  $\log(\pi/(1-\pi))$  [if  $C_P = C_N = 1$ ]

No-pre-Application Tuning

↳ no need to optimize threshold → well  
calibrated model

SE IL MODELLO  
È OTTO INSERITO  
IL DATA SONO  
LONDRA, IL  
TASSO DI CUE  
IL 70% DI  
PERCENTUALE  
X QUALE  
NORDICO

DIFFERENTI  
E IL DECORRERE  
È UGUALE ALLA  
MISURA, IL  
MODELLO È  
BELL'ADATTATO

→ if we can't  
achieve perfect  
calibration in the  
initial training

we apply post-processing  
techniques to improve → learn a mapping  
from the original scores  
to new scores

more like genuine log-likelihood  
ratios or calibrated  
probabilities

$f(s) = f(s)$   
transforming score

Overfitting to Validation Set

↳ validation data used to  
optimize multiple hyperparameters

data no truly unseen → overconfident

changing application contexts

optimal threshold depend  
heavily on specific application → prior probabilities and  
error cost.

we need to recompute minDCF ... (repeat the process)

SE X È SPERATO CONSIDERATO  
IN LOGISTICA REGRESSIONE WEIGHTED  
NON APPLICANDO X UN APPLICATION CHE È SIMILE A QUELLA  
TUTTI I TIPI DI APP SE È POSSIBILE SIMILE LI  
COSTI FAVOREVOLI, SE C'È UNA SIMILE LI  
IL MODELLO CUE HA STATO, COMUNQUE SE È SIMILE  
RISULTATO, COMUNQUE SEMPRE IL TASSO

## Information Preservation

preserve relative ranking of the score

score A > score B

## Meaningful Output

raw score into calibrated posterior probability  $p(c|s)$  → then we can convert into HR easily  
the transformed scores should → log likelihood ratios → transform raw scores directly into calibrated log likelihood ratio

## Monotonic function

↳ calibration doesn't destroy the discriminative information of the original classifier

$\pi_i \in \text{BIN}(\text{UNIFORM}) \text{ OR } \text{CONDENSATI } \pi_i$

$$\pi = \alpha s$$

$$S_2(x) = S_1(x) - 10$$

$x$  CONF è IMPLEMENTAZIONE

$\pi_i$  è UNIFORME

TA TUTTI GLI SCORE

$S_1(x)$  HA LO STESSO SHAPe

MA i COEFFICIENTI A -10, QUINDI

NON HA LA STESSA CATEGORIALE

CONSIDERANDO IN MANIERA CORRETTA

TUTTI I PUNTI DEGLI CLASSI POSSONO

SBALLONO I RESTANTI SONO CLASSE BLU.

L'OPTIMAL THRESHOLD NON È PIÙ IL TRONCO A -10,  $t=0$  WILL DO

very bad!

SO POSSIAMO HAVING UN GOOD ONE HAVING THE SAME SCORE,  $\pi_i^{\text{cal}} = S_2(x) + 10$  così

È BEEN CALIBRATED.

WE SHOULD HAVE UNIFORM DISTRIBUTION OF THE SCORE, ANALOGO ALLA

SI CONSIDERA IL TRANSFORM IN MANIERA TALE DA

NORMALIZZARLO

IN THIS CASE  $\pi_i^{\text{cal}} = S_2(x) + \beta$  / SE SO  
BROUARD AUF IN  
UNIFORM DISTRIBUTION  
NO I MANTIANI  
GIUSTI

## Practical Implementation

How we do this?

## Parametric Calibration Through Score Modeling

Score well calibrated are score that can be thresholded in a good way, this means that if we have a  $\pi=0.5$  and a model that gives us well calibrated score

$S_2(x)$  PIÙ  $\pi_i^{\text{cal}}$  PIÙ UNIFORME  
PIÙ PROBABILITÀ DI UNO APPROPRIATO

We can model the much simpler one-dimensional space of classifier score

We treat score as disjoint and we model their distributions separately per each class

Calibration: prior-weighted logistic regression

We consider the non-calibrated scores as 1-D feature vectors

We assume an affine mapping from non-calibrated scores to calibrated scores

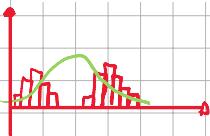
$$f(s) = \alpha s + \gamma$$

Since  $f(s)$  should produce well-calibrated LLRs,  $f(s)$  can be interpreted as the log-likelihood ratio for the two class hypotheses

$$f(s) = \log \frac{f_{S|C}(s|H_T)}{f_{S|C}(s|H_F)} = \alpha s + \gamma$$

The class posterior probabilities for prior  $\bar{\pi}$  correspond to

$$\log \frac{P(C = H_T|s)}{P(C = H_F|s)} = \alpha s + \gamma + \log \frac{\bar{\pi}}{1 - \bar{\pi}} = \alpha s + \beta$$



SE I NUOPELLI GENERATIVI, QUESTI POSSONO AVERE SCORI NON CALIBRATI WHEN THERE IS DISAGREEEMENT BETWEEN THE TRAINING AND THE APPLICATION DATA O NON PIASTRANO DENTI I OGNI OS

BUON SE ALLENAMO IL MODELLO SUL NOSTRO SAMPLE, O VERSO SONO SOLAMENTE UNI DI DIMENSIONE, QUINDI AVEREMO UN BUON FIT E' DIFFICILE Non è core il caso UNIDIMENSIONALE. Non abbiamo nemmeno un numero di sample sufficienti a usare

abbinare  
in realtà

LE PROBLEMI A DESTRA LE FEATURE SE SÌ I NUOPELLI CAUSANTI FANNO CICARE

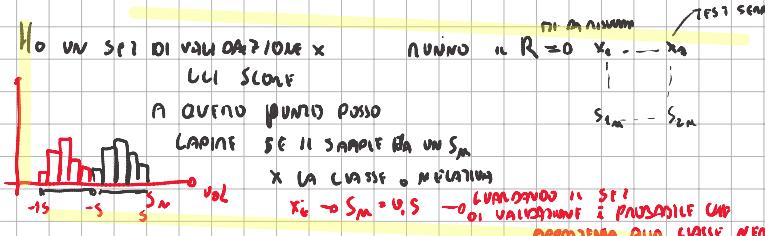


Cose fare? Se non consideriamo la possibilità di mismatch, e consideriamo 2 dimensional data ottenendo degli score che sì well calibrated

$$\text{IL NOSTRO } R = S_m = \log \frac{f(x|H_T)}{f(x|H_F)} \quad (S_m) \text{ NON È UNO SCORO QUINDI}$$

EFFETTIVAMENTE NON SO QUANDO E' CORE ATTINGERE UN DETERMINATO VALORE AD UNA CLASSE POSITIVA O NEGATIVA

$S_m \geq t$ ? Trasformo il problema nel trovare una probabilità per cui so ovvero i sample mappe furono una classe negativa o una positiva.



I SOLO I QUINDI OVRUTI SI

PROBABILMENTE SONO SÌ PIÙ LUNGHE ANNO  
WELL CALIBRATED SCORE (LEI PUOTONO SÌ ESSERE IN UN SETTORE IN MODO TALE CHE NON HOVANO)

X I NUOPELLI GENERATIVI, QUESTI POSSONO AVERE SCORI NON CALIBRATI WHEN THERE IS DISAGREEEMENT BETWEEN THE TRAINING AND THE APPLICATION DATA O NON PIASTRANO DENTI I OGNI OS

INVECE DI USARE IL THRESHOLD, USO IL SET DI VALIDAZIONE (IL THRESHOLD NON SO CONVERGONO LA MIGLIOR VALIDAZIONE MA POSSO I MIGLIOR SCORE M NON SO IN QUESTA SITUAZIONE) SE HO UN SET DI VALIDAZIONE SO IN QUESTA SITUAZIONE MI TROVO,



Noi NON GUARDIAMO L'ISTOGRAMMA, MA GUARDIAMO

$f(S_m|H_T)$  PENSAMO DELL'ISTRADIZIONE SULL'ISTOGRAMMA E NOI TRAVIAMO UN LIKELIHOOD O CLASSE CHE APPARTIENE ALLA CLASSE POSITIVA O NEGATIVA.

LOG  $\frac{f(S_m|H_T)}{f(S_m|H_F)}$  È UNA LOGO CHE DESCRIVE THE DISTRIBUTION OF THE SCORE. SO I HAVERE TO CALCULATE THIS NUMBER ON TO SCORE USARSI A DATASET AND CORRECT THE LIKELIHOOD TO THEN COMPUTE THE CLASSIFICATION.

In this case the model is well calibrated now

COSÌ FUNZIONA TUTTO ORA???

Training sample

$$x_1 - x_m$$

da cui costituisco  
una mia classe  
di dati

$$S_m$$

Validation set  
calcolando UNICO UNICO  
 $x_1^v - x_n^v$  score on the validation  
set

$$= \frac{N(S_m|x_1^v, \mu_1)}{N(S_m|x_1^v, \mu_2)}$$

$S_m(x)$

$$es R \rightarrow N(x|\mu_1, \Sigma_1), N(x|\mu_2, \Sigma_2)$$

CLASSIFICATION

$$x_v: S_m(x_v) = \log \frac{N(x_v|\mu_1, \Sigma_1)}{N(x_v|\mu_2, \Sigma_2)} = \text{Scal}(x_v) = \log \frac{N(S_m(x_v)|\mu_1, \Sigma_1)}{N(S_m(x_v)|\mu_2, \Sigma_2)}$$



How well this model performs depends on how well I predicted the distribution of my score

La nostra trasformazione prima tutta, la mia set di test si trova non calibrata se pur di averli tutti in alto avranno tutti alti minuti

Applicazione all' SVM

$$x_1 - x_m \text{ Training}$$

$$\left\{ \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - z_i(w^T x_i + b)) \right\}$$

Val set

$$x_1^v - x_n^v$$

$$w^T x_1^v - w^T x_n^v$$

3

w<sub>i,b</sub>

$$S_m(x_t) = w^T x_t + b$$

FONDO CAUSSIANE  
 + 10 SVI NUOVI  
 score  
 $S_{cal} = \log f(S_m(x_t), v_p)$

$$\log f(S_m(x_t), v_p)$$

So I HAVE TO FIND A GOOD MODEL FOR THE SCORE. AND USUALLY A GAUSSIAN MODEL IS NOT ALWAYS THE RIGHT DECISION, I CAN USE OTHER MODELS THAT I HAVE SEEN DURING THE COURSE.

For example if I assume that the score distributions have the same covariance

$$\frac{\log f(S_m(x_t)) / \mu_p, v_p}{f(S_m(x_t)) / \mu_p, v_p} = \alpha S_m(x_t) + \beta$$

LINEAR BOWMPTON

affine mapping from non-calibrated to calibrated, usually a linear transformation is good enough.

We can use a logistic regression using affine function

$$\sigma(S_{cal}(x_t)) = \frac{P(c = H_t | x_t)}{P(c = T_t | x_t)}$$

on our  $x$  obtaining  $S_{un}(x_t) = S_{cal}(x_t) - \log \frac{\pi}{1-\pi}$

A SUMMATA SE NON VOLGONO USARE LA SOTTRAZIONE CON LA LOGIT DOPO PONIAMO APPLICARE LA FORMULA WEIGHTED

The calibration transformation corresponds to the transformation of a prior-weighted logistic regression model

Once we have estimated  $\alpha$  and  $\beta$ , the calibrated score  $f(x)$  can be computed as

$$f(x) = \alpha x + \gamma = \alpha x + \beta - \log \frac{\pi}{1-\pi}$$

Note that we have to specify a prior  $\pi$ , and we are still effectively optimizing the calibration for a specific application  $\pi$ . However, often this approach will provide good calibration for a wider range of different applications similar to the target one

We can also modify the prior-weighted logistic regression objective to compute directly  $\alpha$  and  $\gamma$ :

$$R(\alpha, \gamma) = \sum_i w_i \log \left( 1 + e^{-z_i(\alpha x_i + \gamma)} \right), \quad w_i = \begin{cases} \hat{y}_i / \hat{\pi} & \text{if } z_i = +1 \\ (1 - \hat{y}_i) / (1 - \hat{\pi}) & \text{if } z_i = -1 \end{cases}$$

IN QUESTO CASO RISULTEREMO  $\hat{\pi}_t$   
 E NON TEMP TO CALIBRAZIONE

IN QUESTO CASO CI POTREBBE ESSERE PROBLEMA DI CALIBRAZIONE?? SE LA CN OVERFITS!! ALLORA HANNO UNA RELAZIONE STRONG BUT IN THIS CASE PUÒ PORTARCI AD UNDERFITTING, IL SCORE È POCO CALIBRATO

MA QUANDO HO EFFETTIVAMENTE OVERFITTING? C'È UN CASO? LA DIMENSIONALITÀ DEI DATI  
 CI SONO FEATURE TANTE QUANTO BISOLITO DI PIÙ PARTEMI X LA STIMA, MA SE LA DIMENSIONE È PIUTTOGLIO  
 NON RISULTANO OVERFITTING  $\Rightarrow$  NON ADDIMO BISOLITO DI UNA RELAZIONE

EXAMPLE:

ASSUMO 2 OLTRENUOVE DATA

• C'È UNO SCORO  
 PIÙ POSITIVO TRAVERSO  
 PIÙ POCHI  
 DATI IN PREDICTION  
 PIÙ FRAGILE NELLE  
 PUÒ ESSERE  
 PIÙ COMPLICATO

G(SUN)

Training set

$$x_1, \dots, x_n$$

$$w, b$$

$$S_m(x) = w^T x + b$$

Val set

$$x'_1, \dots, x'_n$$

$$S_m(x'_1) - S_m(x'_n) \rightarrow \text{diff}$$

SAMPLING

1-0  
 SV (verso) APPROPRIATO  
 DATA POINT WEIGHTED =  
 LOUNGIST REGRESSION

$$\alpha^*, \beta^* = \arg \min_{\alpha, \beta} \sum_{i=1}^n \sum_{j=1}^n \log \left( 1 + e^{-z_j(\alpha x_j + \beta)} \right) \rightarrow \frac{\partial}{\partial \alpha} \sum_{i=1}^n \sum_{j=1}^n \log \left( 1 + e^{-z_j(\alpha x_j + \beta)} \right)$$

$$S_m = S_m(x')$$

$\alpha^*$   $\beta^*$   
 ALLENNO IL MODELLO

$$x_t \rightarrow S_m(x_t) \rightarrow S_{post} \rightarrow \alpha^* S_m(x_t) + \beta^* \rightarrow S_{cal}(x_t) = \alpha^* S_m(x_t) + \beta^* - \log \frac{\pi}{1-\pi}$$

se alleno questo non  
 è il modello di  
 TOLLING  
 $\frac{\pi}{1-\pi}$

MA È LA STESSA COSA

OULAVANTE SP HO UN CASO così che molti CDF POSSONO ESSERE CONFERMANTE SEPARATI, ALLORA HANNO OVERFITTING, (AVREMO DISORDINE DELLA REGULARIZZAZIONE) MA SONO GIÀ NON ACCADE. (AVREMO DISORDINE PIÙ OGNI)

Nota: IL MODELLO DAQI SCORE WILL CALCULATE = UNA VARIABILE DIFERENTI IN BASE ALA COMPLESSITÀ DELLA MAPPING FUNCTION.

Se ho applicato dove la misura è 0,1 il modello lavorerà bene anche per  
 APP. 0,01, se è molto diverso però dovrà allentare un altro modello.

SE L'APPARTENENZA DELL'  
 DISTINZIONE È DUREVA  
 ALLORA VA BENE X  
 PIÙ APPROPRIATO

GIÀ NON AGRADE X I MODULI CENTRALI PIÙ

## Score models

- Require assumptions on the calibration transformation (e.g. affine mapping) or on the distribution of class scores
- Estimated over a training set, but allow for extrapolation outside of training score ranges
- Typically, fast to evaluate
- May provide good fit only for a small range of operating points

### Example: Prior-weighted logistic regression

$$v_{\text{univ}} \quad v_{\text{F}} = \alpha s + \beta$$

### Calibration: prior-weighted logistic regression

We consider the non-calibrated scores as 1-D feature vectors

We assume an affine mapping from non-calibrated scores to calibrated scores

$$f(s) = \alpha s + \gamma$$

Since  $f(s)$  should produce well-calibrated LLRs,  $f(s)$  can be interpreted as the log-likelihood ratio for the two class hypotheses

$$f(s) = \log \frac{f_{S|C}(s|\mathcal{H}_T)}{f_{S|C}(s|\mathcal{H}_F)} = \alpha s + \gamma$$

The class posterior probabilities for prior  $\bar{\pi}$  correspond to

$$\log \frac{P(C = \mathcal{H}_T | s)}{P(C = \mathcal{H}_F | s)} = \alpha s + \gamma + \log \frac{\bar{\pi}}{1 - \bar{\pi}} = \alpha s + \beta$$

achieving low calibration error

↳ Atention characteristics of the dists  
useful for calibration

↳ calibration set should  
be similar to evaluation set

eval and use unseparated  
to avoid overfitting

- Miscalibration due to non-probabilistic scores, or to overfitting or underfitting models, but evaluation and training populations are similar: the calibration set can be extracted from the training set material

In this case, calibration allows recovering a probabilistic interpretation of the scores, and, since we are training on 1-D data, the risk of overfitting is drastically reduced (and typically we don't need to add regularization to the calibration model objective function)

we use this

We can employ the (non-regularized) prior-weighted logistic regression model to learn the model parameters  $\alpha, \beta$  from a set of training scores (we will refer to the set of samples employed to estimate the calibration parameters as **calibration training set** in the following)

In practice, we are treating scores as if they were 1-D feature vectors

As for model training and validation set, also the calibration set should be an independent dataset that does not overlap with either the model training nor the validation set (we will see later how to effectively split the data)

The calibration transformation corresponds to the transformation of a prior-weighted logistic regression model

Once we have estimated  $\alpha$  and  $\beta$ , the calibrated score  $f(s)$  can be computed as

$$f(s) = \alpha s + \gamma = \alpha s + \beta - \log \frac{\bar{\pi}}{1 - \bar{\pi}}$$

Note that we have to specify a prior  $\bar{\pi}$ , and we are still effectively optimizing the calibration for a specific application  $\bar{\pi}$ . However, often this approach will provide good calibration for a wider range of different applications similar to the target one

We can also modify the prior-weighted logistic regression objective to compute directly  $\alpha$  and  $\gamma$ :

$$R(\alpha, \gamma) = \sum_i w_i \log \left( 1 + e^{-\alpha(z_i + \gamma + \log \frac{\bar{\pi}}{1 - \bar{\pi}})} \right), \quad w_i = \begin{cases} \bar{\pi}/n_T & \text{if } z_i = +1 \\ (1 - \bar{\pi})/n_F & \text{if } z_i = -1 \end{cases}$$

Alternatively to logistic regression, we can employ generative models to estimate calibrated LLRs

We model the class-conditional **score** distribution for the two classes

$$S|\mathcal{H}_T, \quad S|\mathcal{H}_F$$

For example, we can employ a Gaussian model to estimate the two densities

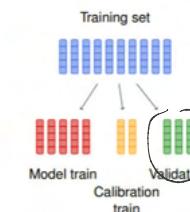
$$f_{S|\mathcal{H}_T}(s) = \mathcal{N}(s|\mu_T, \nu_T), \quad f_{S|\mathcal{H}_F}(s) = \mathcal{N}(s|\mu_F, \nu_F)$$

and the calibration transformation is given by the LLR for the score model

$$f_{\text{cal}}(s) = \log \frac{f_{S|\mathcal{H}_T}(s)}{f_{S|\mathcal{H}_F}(s)}$$

Tied models  $v_T = v_F = v$  are usually employed, as they result in monotonic transformations

103



use ONLY  
set DCF:  
The metric  
that we will  
actually use

Validation set

for evaluation

we use set DCF no min DCF

the score have been  
calibrated.

total → 3 set Framework for  
calibration

Training set → for learning

Calibration set → used for calibrate the raw score  
true posterior probabilities

- Miscalibration due to non-probabilistic scores, or to overfitting or underfitting models, but evaluation and training populations are similar: the calibration set can be extracted from the training set material

We need to collect data similar to the application population

The amount of required matching data, however, is usually small compared to the amount required to train a complete classifier

Some models (e.g. Gaussian) can be extended to exploit unlabeled samples (unsupervised or semi-supervised training), which are less expensive to acquire

How it work? → In lab → We have solved the score compute on the validation set → subdivide the validation set into new partitions

train the classifier → then we apply it to calibration → feature for partition to generate score  
apply calibration → the apply the technique  
model to validation partition → use calibrated score method to generate the score end

are simple one for calibration → validation transformation

the other for validation evaluation

big problem

progressive reduction in dataset size available  
underfitting/overfitting

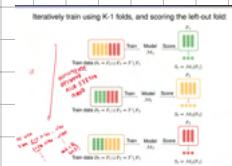
each fold maintains the statistical characteristic of the original dataset  
same prior  
training and evaluation cycle works with representative subsets that reflect over class distribution

For large dataset random sampling is sufficient

start with dividing the dataset into K approximately equal partitions → folds

perfect equality is not required

1,000 samples into three folds of sizes 333, 333, and 334



with data size increase model train on overlapping subset of data should coverage to similar solution

with can use CV for hyperparameter tuning → We use the same procedure for if  $\text{svr} = \epsilon$ ,  $C=10$  → and based on the DCF value we can choose

Pool together the scores of the different folds and evaluate the desired metric

Use the chosen metric to perform model selection / hyperparameter tuning for each model

NOTE (1): all the  $K$  models  $M_1, \dots, M_K$  must be trained with the same set-up (i.e., same pre-processing, same values for the hyperparameters, ...)

each model that we train can make different prediction

aggregated scores may not accurately represent what any single model would produce on the same evaluation data

Solution → Cross Validation  
virtually use all the data for both training and validation

no more limitation data can serve for only one purpose

advantage

eliminate trade off → training and evaluation benefit → independence is preserved

maximum possible data

We can also compare model SVM vs LR

first fold that has only that sample has a very low DCF  
because it is capable of data east to classify

I CANT USE DCF to choose  
ALWAYS I WILL NOT CHOOSE ANY OF THE REST FOLD

to have the best accuracy when we need to reduce the dimension of the folds

we maximize the similarity between cross-validation models

leave one out is the best sample per fold

validation fold where relatively small, compared to training sets

additional data should have minimum impact on the model characteristics (hyperparameter computed with cross-validation)

Solution: we train one more model over the whole training set, using the method we selected and the hyperparameters we estimated in the previous step



The model  $M$  will then be used to compute the scores for the application data

For this reason we aggregate the OCR of each fold

and finally for the same reason

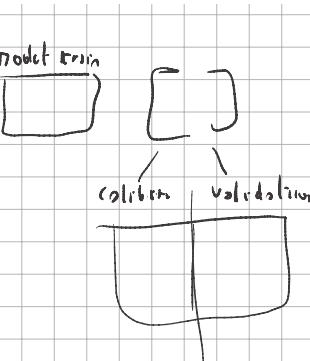
how we compute the training parameter?  $w$  and  $b$

Unfortunately, leave-one-out may require training an excessively large number of models (we need to train  $N$  times a model over  $N - 1$  samples)

We choose  $K$  as large as we can afford

- Large values of  $K$ : more robust analysis, more time required
- Small values of  $K$ : less robust analysis, but faster

→ In laboratory we will use single-folds and  $K$ -folds



## Classifier Fusion → combining the prediction of multiple classifiers

CM make some mistakes in some cases

SVN in other cases

This complementary error pattern creates opportunity for improving

How we do it?

Majority voting

What happens if we have only two classifiers?  
We need to consider also the confidence

$C_1$	$H_1$	$C_1 S = 0.1$	$\pi_1 = 0.5$
$C_2$	$H_2$	$C_2 S = 0.05$	
$C_3$	$H_3$	$C_3 S = -0.5$	

I should rely on the most confident

- QUESTION È PIÙ CONFIDANTE  
ANCHE SE HA L'ERRORE  
P1 DEDICO ALLE  
4 CLASSI VERA E NUDA

We have 3 classifier

$$\begin{array}{c} n_1 \quad n_2 \quad n_3 \\ s_1 \quad s_2 \quad s_3 \end{array}$$

IL PIÙ BUONI?

come unendo tutti  
score?

it's a correct approach  
when classifier provide →  
independent information

(the model)  
need to be independent  
Es: no produce a score looking  
at the first half of the features →  
→  $n_1$  at the second half

what we will do →

We can consider weighting the contribution of the different systems. Given scores  $s_1 \dots s_m$  for a test sample  $x$ , we can compute a "fused" score as

$$s_{\text{fused}} = \alpha_1 s_1 + \alpha_2 s_2 + \alpha_3 s_3 + \dots + \alpha_m s_m + \gamma$$

We have the problem of estimating the weights  $\alpha_i$  and the bias term  $\gamma$

If we stack the scores and the weights as

$$s = \begin{bmatrix} s_1 \\ \vdots \\ s_m \end{bmatrix}, \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix}$$

the fused score can be represented as

$$s_{\text{fused}} = \alpha^T s + \gamma$$

USANDO UN MODELLO X DARE UN PESO  
AD Ogni CLASSIFICATORE

!!

COSÌ CALCOLANDO  $\gamma$ , I PARAMETRI NEL  
QUALE NUOVO?

$$\frac{n_1}{n} = p$$

$$\begin{aligned} n_1 -> S = 10 \\ n_2 -> S = 10 \end{aligned}$$

IF I SUM  $n_1 + n_2 = 20$  → THIS SCORE  
WILL NOT BE MUCH CORRELATED  
ANTHONY

IN THIS CASE WE SHOULD USE THE  
AVG APPROX

problem if  
classifiers are  
correlated  
we can't use fusion model

$$\begin{aligned} & x \\ & \left. \begin{array}{l} x^+ \\ x^- \end{array} \right\} n_1 \quad \left. \begin{array}{l} x^+ \\ x^- \end{array} \right\} n_2 \\ & \frac{f(x|n_1)}{f(x|n_2)} = \frac{f(x^+|n_1) f(x^-|n_1)}{f(x^+|n_2) f(x^-|n_2)} \\ & = \log \frac{f(x^+|n_1)}{f(x^+|n_2)}, \dots \rightarrow s_{n_1} + s_{n_2} \end{aligned}$$

