

## ESTIMATING THE GENERALIZATION ERROR

- For a fixed model, we are interested in the Generalization Error (GE):  $GE(\hat{f}, L) := \mathbb{E} [L(y, \hat{f}(\mathbf{x}))]$ , i.e. the expected error the model makes for data  $(\mathbf{x}, y) \sim \mathbb{P}_{xy}$ .
- We need an estimator for the GE with  $m$  test observations:

$$\widehat{GE}(\hat{f}, L) := \frac{1}{m} \sum_{(\mathbf{x}, y)} [L(y, \hat{f}(\mathbf{x}))]$$

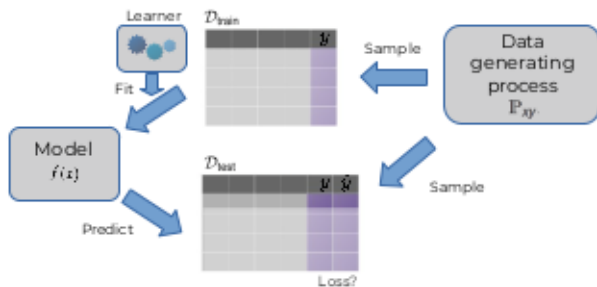
- However, if  $(\mathbf{x}, y) \in \mathcal{D}_{\text{train}}$ ,  $\widehat{GE}(\hat{f}, L)$  will be biased via overfitting the training data.
- Thus, we estimate the GE using unseen data  $(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}$ :

$$\widehat{GE}(\hat{f}, L) := \frac{1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} [L(y, \hat{f}(\mathbf{x}))]$$



## ESTIMATING THE GENERALIZATION ERROR / 2

- Usually, we have no access to new **unseen** data.
- Thus, we divide our data set manually into  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$ .
- This process is depicted below.



## METRICS FOR CLASSIFICATION / 2

For hard-label classification, the confusion matrix is a useful representation:

		True Class $y$	
		+	-
Pred. $\hat{y}$	+	True Positive (TP)	False Positive (FP)
	-	False Negative (FN)	True Negative (TN)



From this matrix a variety of evaluation metrics, including precision and recall, can be computed.

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

## ESTIMATING THE GENERALIZATION ERROR (BETTER)

While

$$\widehat{\text{GE}}(\hat{f}, L) := \frac{1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} [L(y, \hat{f}(\mathbf{x}))]$$

will be unbiased, with a small  $m$  it will suffer from high variance. We have two options to decrease the variance:

- Increase  $m$ .
- Compute  $\widehat{\text{GE}}(\hat{f}, L)$  for multiple test sets and aggregate them.

With a finite amount of data, increasing  $m$  would mean to decrease the size of the training data. Thus, we focus on using multiple ( $B$ ) test sets:

$$\mathcal{J} = ((J_{\text{train},1}, J_{\text{test},1}), \dots, (J_{\text{train},B}, J_{\text{test},B})) .$$

where we compute  $\widehat{\text{GE}}(\hat{f}, L)$  for each set and aggregate the estimates. These  $B$  sets are generated through **resampling**.

