



香港城市大學  
City University of Hong Kong

專業 創新 胸懷全球  
Professional • Creative  
For The World

# Evaluation: The Problem of Overfitting

CS5483 Data Warehousing and Data Mining

2

## Guess the value of $y$

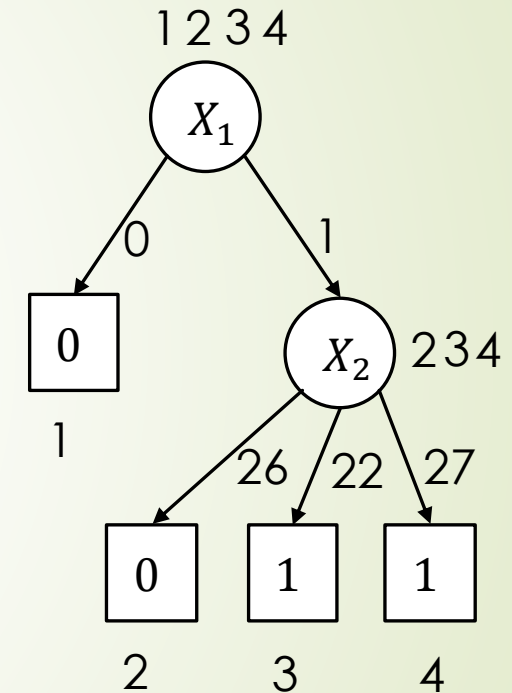
	$X_1$	$X_2$	$Y$
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1
	1	22	$y$

►  $y = \underline{\hspace{1cm}}$  because                     .

3

# Guess the value of $y$ for any $(x_1, x_2)$

	$X_1$	$X_2$	$Y$
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1
	$x_1$	$x_2$	$y$

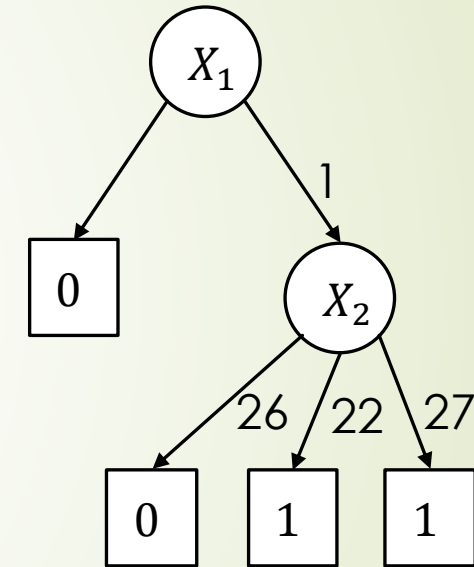


- D\_\_\_\_\_t\_\_\_\_\_.
- Is it making good decisions? It **f**\_\_\_\_ the data **p**\_\_\_\_\_.
- Is it simple? \_\_\_\_\_

4

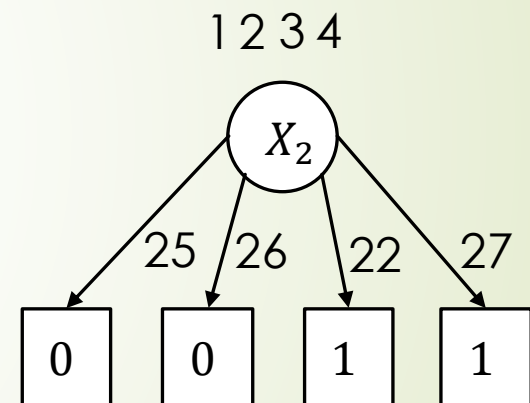
Guess the value of  $y$  for any  $(x_1, x_2)$

	$X_1$	$X_2$	$Y$
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1
	$x_1$	$x_2$	$y$



# Guess the value of $y$ for any $(x_1, x_2)$

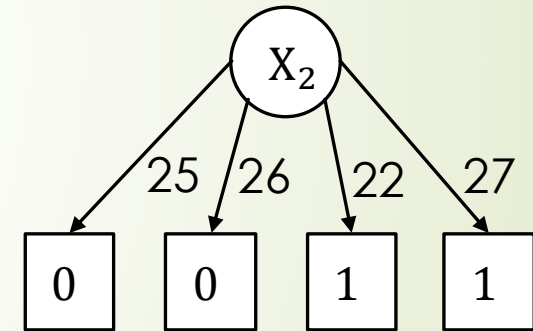
	$X_1$	$X_2$	$Y$
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1
	$x_1$	$x_2$	$y$



- Can have a smaller tree because  $X_2$  completely \_\_\_\_\_  $Y$ .
- Smaller tree means faster computation and lower storage.
- Any other benefit? What if  $(x_1, x_2) = (0, 22)$ ?  $y = \underline{\hspace{1cm}}$  (generalize to u\_\_\_\_\_ data.)

# Does the classifier generalize well?

	$X_1$	$X_2$	$Y$
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1
	$x_1$	$x_2$	$y$

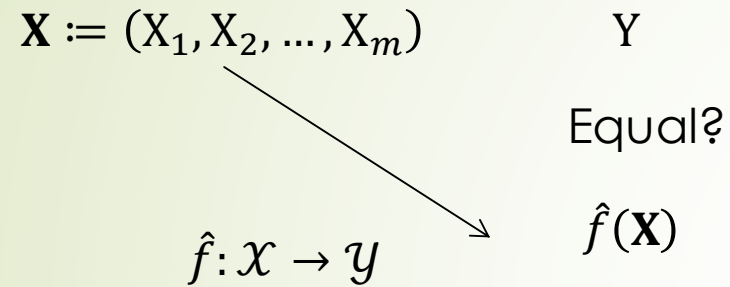


- $Y$ : Diagnosis of certain disease.
- $X_1$ : Result of a medical test.
- $X_2$ : Temperature at which the test is conducted.
- If  $(x_1, x_2) = (0, 22)$ , should  $y = 1$ ? Yes/No, because \_\_\_\_\_.

# Overfitting

- Learning patterns that
  - **fits** the data well, but
  - not **g**\_\_\_\_\_ well to new cases (future data).
- The challenges:
  1. How to estimate the actual performance on unseen data, not the **overly-optimistic** performance on fitted data?
  2. How to learn the desired knowledge, i.e., patterns that **generalize** well to unseen data?
- What causes overfitting?

# Naive Formulation



- **Input feature** vector  $\mathbf{X}$  from **feature space**  $\mathcal{X}$ .
- **Target**  $Y$  from a discrete set  $\mathcal{Y}$ :
  - B\_\_\_\_\_ classification:  $|\mathcal{Y}| = 2$
  - M\_\_\_\_\_ classification:  $|\mathcal{Y}| > 2$
- Obtain a **classifier**  $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$  to predict  $Y$  from  $\mathbf{X}$ .
- What is a good classifier?



# Naive Formulation

$\mathbf{X} := (X_1, X_2, \dots, X_m)$

$Y$

$\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$

$\hat{f}(\mathbf{X})$

$$\hat{f} \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} P[Y \neq f(\mathbf{X})]$$

is a **classifier** that

- Minimize  $\Pr[Y \neq f(\mathbf{X})]$ , **e**\_\_\_\_\_ probability.
- Maximize  $\Pr[Y = f(\mathbf{X})]$ , **a**\_\_\_\_\_.

- More generally, choose  $\hat{f} = f_w$  where  $w$  is a solution to

$$\min_{w \in \mathcal{W}} E \left[ \overbrace{\underbrace{L(Y, f_w(\mathbf{X}))}_{R(w)}}^{\mathbf{1}(Y \neq f_w(\mathbf{X}))} \right]$$

- $L: \mathcal{Y}^2 \rightarrow \mathbb{R}$  can be any **l**\_\_\_\_ function such as the 0-1 loss  $(y, \hat{y}) \mapsto \mathbf{1}(y \neq \hat{y})$ .
- $R: \mathcal{W} \rightarrow \mathbb{R}$  is the corresponding **r**\_\_\_\_ functional over a (compact) hypothesis space  $\mathcal{W}$  for any class  $\{f_w | w \in \mathcal{W}\}$  of functions.

# Naive Formulation

$$\mathbf{X} := (X_1, X_2, \dots, X_m) \quad Y$$

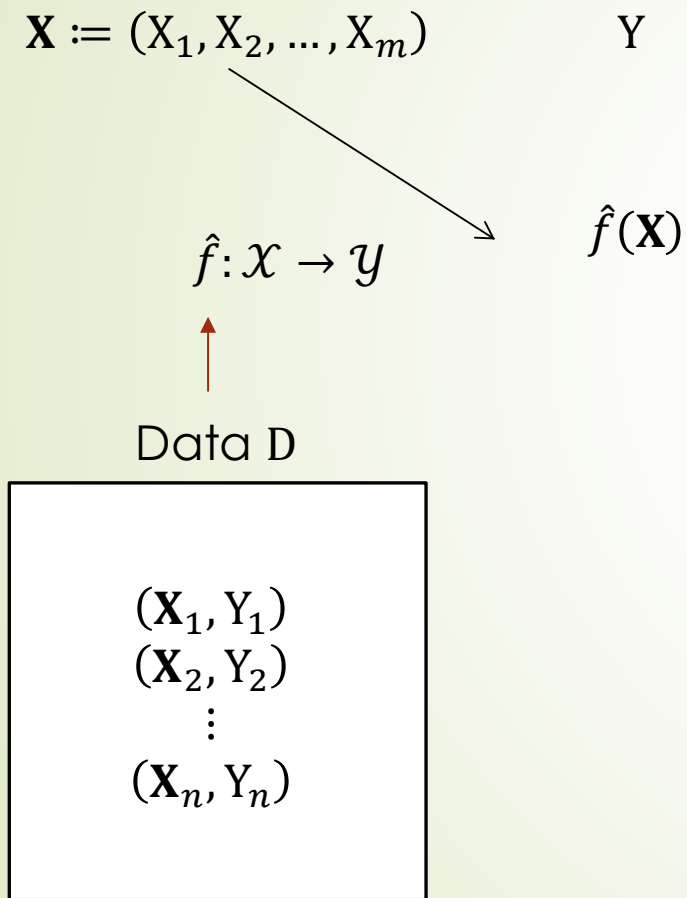
$w \in \mathcal{W} \searrow$

$$f_w(\mathbf{X})$$

$$\min_{w \in \mathcal{W}} \overbrace{\int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} L(y, f_w(x)) dP_{\mathbf{X}\mathbf{Y}}(x, y)}^{R(w)}$$

- Optimal? Yes/No because \_\_\_\_\_.
- Cannot use in practice because the **d**\_\_\_\_\_ is often unknown.

# Empirical Risk Minimization (ERM)



- Obtain/train a **classifier**  $\hat{f}$ 
  - without knowing  $P_{\mathbf{X}Y}$ , but
  - with **data**  $D$  consisting of *i.i.d.* samples of  $(\mathbf{X}, Y)$  and *independent* of  $(\mathbf{X}, Y)$ .
- Fundamental questions:
  - What is a good classifier?
  - How to train a good classifier?

# Estimate the probability of head



Unknown  
 $p := P[\text{Head}]$

Toss  $n$  times

H, T, H, ...

Estimate  $p$

$\hat{p}$

- From the outcomes of  $n$  independent coin tosses, how to estimate the probability  $p$  of the coin coming up head?

$\hat{p} := \underline{\hspace{1cm}}$  in terms of the number  $N_H$  of heads in  $n$  coin tosses.

- How good is the estimate?

# Estimate expectation by sample average

- Given i.i.d.  $n$ -sample of  $Z$

$$Z^n := (Z_i | i \in [n]) = (Z_1, \dots, Z_n),$$

estimate the expectation  $E[Z]$  by

$$\bar{Z} := \frac{1}{n} \sum_{i=1}^n Z_i, \quad (\text{sample average})$$

- The estimate is

- Unbiased, i.e.,

$$E[\bar{Z}] = E[Z]$$

- Consistent, i.e.,

$$\lim_{n \rightarrow \infty} \bar{Z} = E[Z]$$

- $\hat{p} := \frac{1}{n} \sum_{i=1}^n Z_i$  and  $E[Z] = p$  by defining the indicator random variable

$$Z := \begin{cases} 1, & \text{—} \\ 0, & \text{—} \end{cases}$$

# Empirical Risk Minimization (ERM)

- Estimate risk from data

- **Empirical risk:**

$$\hat{R}(w) := \frac{1}{|D|} \sum_{(x,y) \in D} L(y, f_w(x)) = \frac{1}{n} \sum_{i=1}^n L(Y_i, f_w(\mathbf{X}_i))$$

- **Empirical error probability** with 0-1 loss:

$$\hat{R}(w) = \frac{1}{|D|} \sum_{(x,y) \in D} \mathbb{1}(y \neq f_w(x)) = \frac{|\{i \in [n] | Y_i \neq f_w(\mathbf{X}_i)\}|}{n}$$

- Choose the best classifier  $f_{\hat{w}}$  where

$$\hat{w} \in \arg \min_{w \in \mathcal{W}} \hat{R}(w)$$

- How good is  $\hat{f}$ ? \_\_\_\_\_

# Performance estimate

$\hat{R}(\hat{w})$  is a **good** estimate of  $E[L(Y, f_{\hat{w}}(\mathbf{X}))]$ ?

- Is the empirical risk consistent, i.e.,

$$\hat{R}(\hat{w}) = \frac{1}{n} \sum_{i=1}^n L(Y_i, f_{\hat{w}}(\mathbf{X}_i)) \xrightarrow{n \rightarrow \infty} E[L(Y, f_{\hat{w}}(\mathbf{X}))]?$$

\_\_\_\_\_ because  $L(Y_i, f_{\hat{w}}(\mathbf{X}_i))$  \_\_\_\_\_ independent over  $i \in [n]$ .

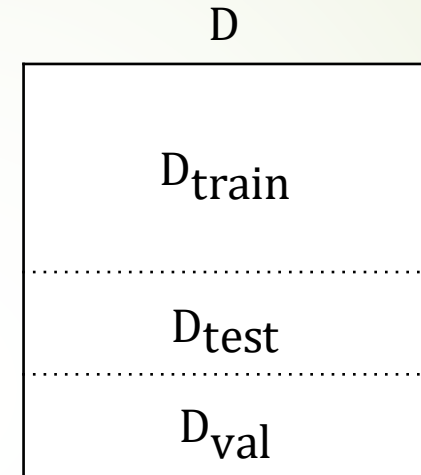
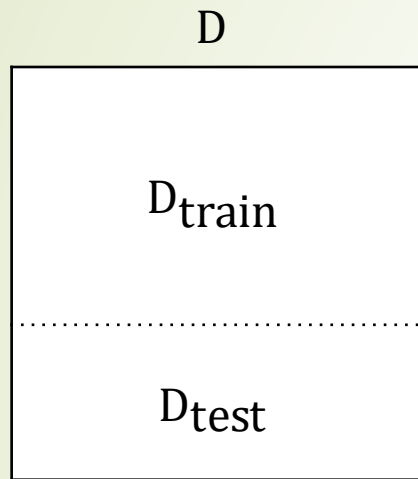
- Is the empirical risk unbiased, i.e.,

$$E[\hat{R}(\hat{w})] = E\left[\frac{1}{n} \sum_{i=1}^n L(Y_i, f_{\hat{w}}(\mathbf{X}_i))\right] = \frac{1}{n} \sum_{i=1}^n E[L(Y_i, f_{\hat{w}}(\mathbf{X}_i))] = E[L(Y, f_{\hat{w}}(\mathbf{X}))]?$$

\_\_\_\_\_ because  $L(Y_i, f_{\hat{w}}(\mathbf{X}_i))$  \_\_\_\_\_ identically distributed as  $E[L(Y, f_{\hat{w}}(\mathbf{X}))]$ .



# Holdout method



- Hold out some data for testing.
  - $D_{\text{train}}$  : **Training set** for constructing the classifier.
  - $D_{\text{test}}$  : **Test set** for testing the classifier.
  - $D_{\text{val}}$  : **Validation set** sometimes for model selection.
- Usually 2: 1 or 2: 1: 1 split. With abundant data, can be 9: 1.



# Empirical risk on a separate test set

- Train a classifier on  $D_{\text{train}}$ , e.g.,

$$\hat{w} \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \frac{1}{|D_{\text{train}}|} \sum_{(x,y) \in D_{\text{train}}} L(y, f(x)).$$

- Compute the empirical risk on  $D_{\text{test}} := ((\mathbf{X}_{s_i}, \mathbf{Y}_{s_i}) | i \in [n'])$

$$\hat{R}(w) := \frac{1}{|D_{\text{test}}|} \sum_{(x,y) \in D_{\text{test}}} \mathbb{1}(y \neq f(x)) = \frac{1}{n'} \sum_{i=1}^n L(\mathbf{Y}_{s_i}, f_w(\mathbf{X}_{s_i}))$$

- which is unbiased, i.e.,

$$E[\hat{R}(\hat{w})] = E\left[\frac{1}{n'} \sum_{i=1}^n L(\mathbf{Y}_{s_i}, f_{\hat{w}}(\mathbf{X}_{s_i}))\right] = \frac{1}{n'} \sum_{i=1}^n E\left[L(\mathbf{Y}_{s_i}, f_{\hat{w}}(\mathbf{X}_{s_i}))\right] = E[L(Y, f_{\hat{w}}(\mathbf{X}))]$$

because  $\hat{w}$  is \_\_\_\_\_ of  $\hat{R}$  as  $D_{\text{train}}$  and  $D_{\text{test}}$  are independent.

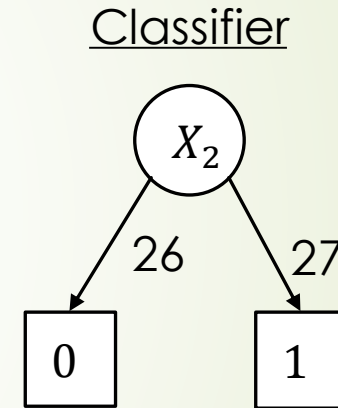
- and consistent, i.e.,

$$\hat{R}(\hat{w}) = \frac{1}{n'} \sum_{i=1}^n L(\mathbf{Y}_{s_i}, f_{\hat{w}}(\mathbf{X}_{s_i})) \xrightarrow{n' \rightarrow \infty} E\left[L(\mathbf{Y}_{s_i}, f_{\hat{w}}(\mathbf{X}_{s_i}))\right]$$

by the uniform law of large number if  $f_w(x)$  satisfies some conditions. (See Theorem 2 [here](#))

# Fails on new values

				<u>Training set</u>				
				$X_1$	$X_2$	$Y$	$\hat{Y}$	err.
1.	0	25	0	2.	1	26	0	0
2.	1	26	0	4.	1	27	1	1
3.	1	22	1	Error rate= _____				
4.	1	27	1	Accuracy= _____%				
				<u>Test set</u>				
				$X_1$	$X_2$	$Y$	$\hat{Y}$	err.
1.	0	25	0	1.	0	25	?	+
3.	1	22	1	3.	1	22	?	+
				Accuracy= _____%				



# Splitting point for numeric attribute

	$X_1$	$X_2$	$Y$	
1.	0	25	0	
2.	1	26	0	
3.	1	22	1	
4.	1	27	1	

Training set

	$X_1$	$X_2$	$Y$	$\hat{Y}$	err.
2.	1	26	0	0	
4.	1	27	1	1	

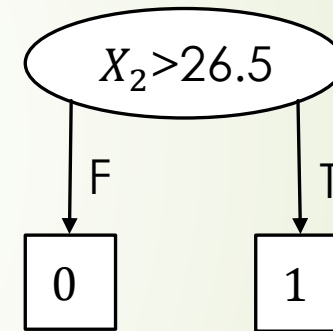
Error rate =  $\frac{0}{2}$   
Accuracy = 100%

Test set

	$X_1$	$X_2$	$Y$	$\hat{Y}$	err.
1.	0	25	0	0	
3.	1	22	1	0	+

Accuracy = \_\_\_\_\_%

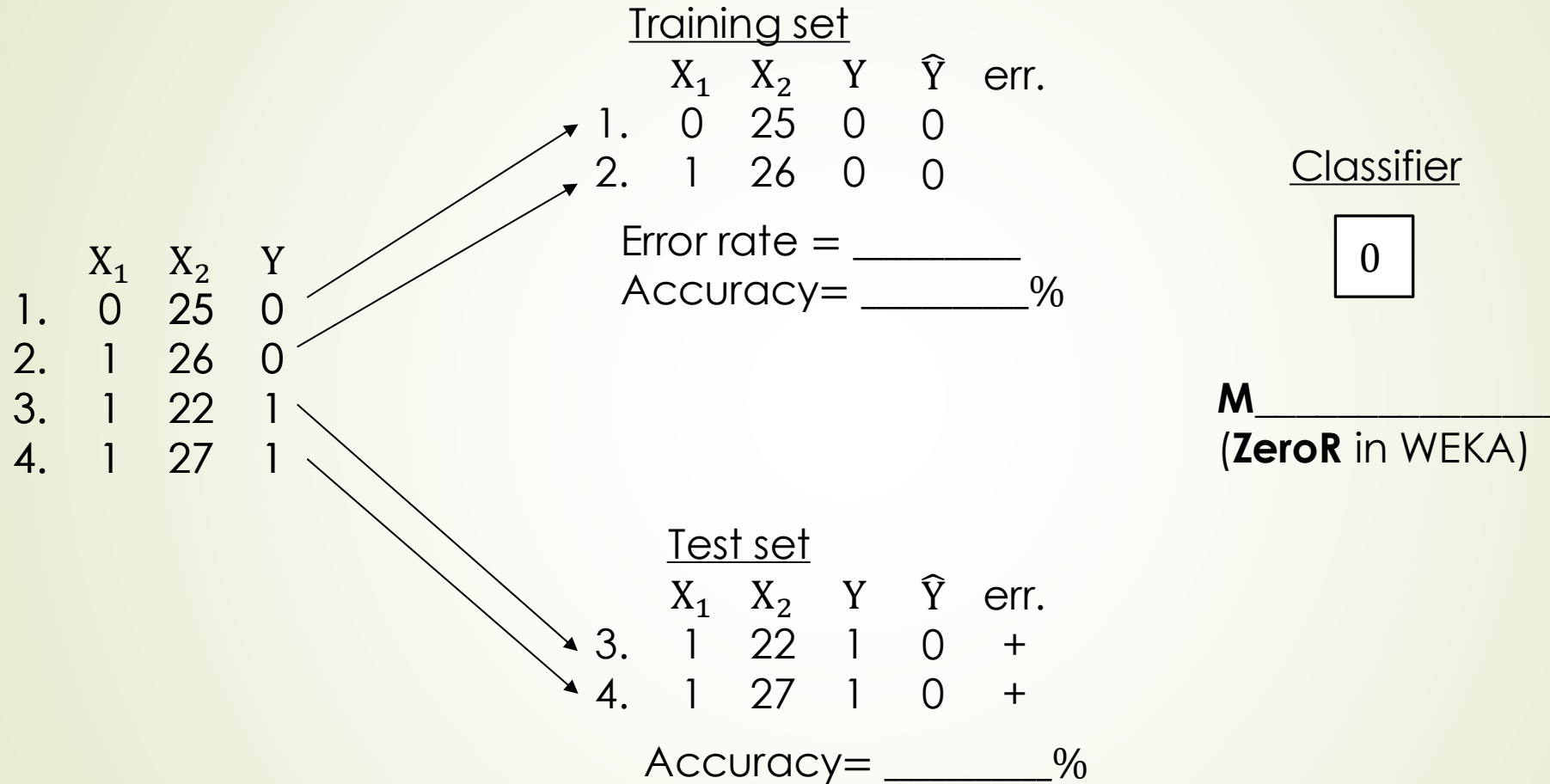
Classifier



# Is 50% accuracy good/bad?

- $Y = 0$  50% of the time.
- A classifier making **a** \_\_\_\_\_ **choices** gives 50% accuracy.

# A different split



- Different splits likely give different classifiers and accuracies.

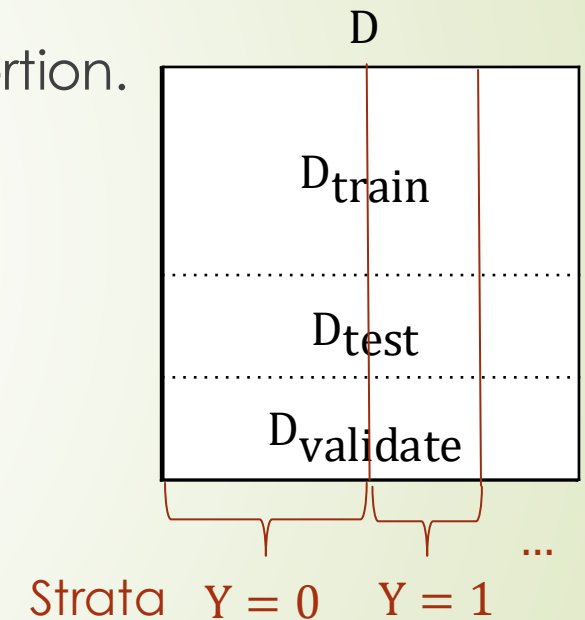
# Class proportion

If the class proportion is not maintained,

- ▶ test set does not reflect the actual performance because  $L(y, f_{\hat{w}}(\mathbf{x}))$  for  $(\mathbf{x}, y) \in D_{\text{test}}$  may not be \_\_\_\_\_ as  $L(Y, f_{\hat{w}}(\mathbf{X}))$ .
- ▶ training/validation set can also mislead the learning process.

# Stratification

- Sample different classes (strata) independently so samples from different classes maintain the class proportion.
- E.g., **Stratified holdout** maintain the proportions of class values in training/test/validation sets.



# How to estimate the typical performance?

- Obtain  $N > 1$  random splits and average the performance.
- **Random s\_\_\_\_\_:**
  - Randomly sampling **without replacement** for training and test sets.
- Training/test set may still be too small.

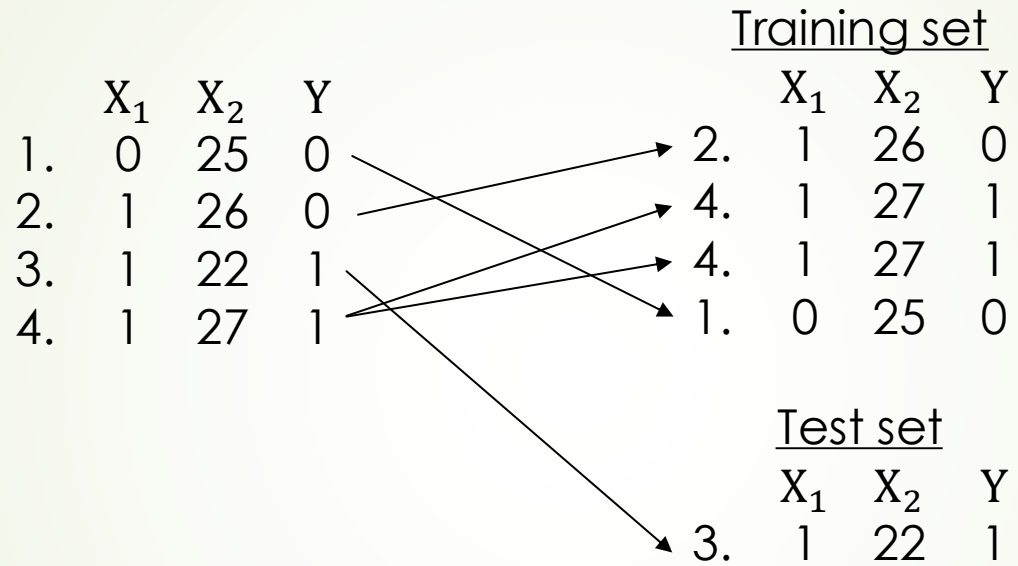


# Can we make training/test sets larger?

➤ **B**\_\_\_\_\_:

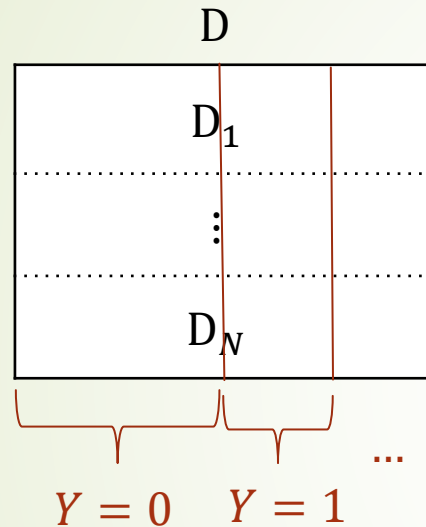
- “pull oneself up by one’s bootstraps”
- Sampling **with replacement** for training set with the same size as  $D$ .
- Remaining **unsampled data** for test set.

## 63.2 bootstrap



- Bootstrap  $n$  training samples.
- Expected portion of data for test is \_\_\_\_\_.

# $N$ -fold stratified cross-validation



Test:  $D_1$ , Training:  $D_2, \dots, D_N$

$\vdots$

Test:  $D_j$ , Training:  $D_1, \dots, D_{j-1}, D_{j+1}, \dots, D_N$

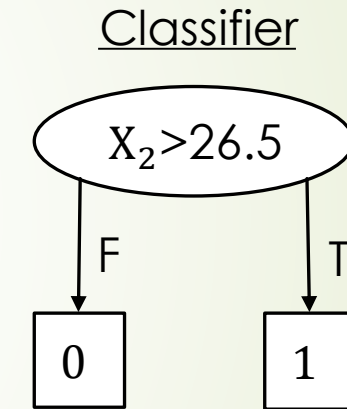
$\vdots$

Test:  $D_N$ , Training:  $D_1, \dots, D_{N-1}$

1. Split data into  $N$  parts/folds equally maintaining the class proportion.
2. Repeatedly take a different fold to test on a model trained on all other parts.
3. Report the average performance on the  $N$  folds.

# Example of 2-fold stratified cross-validation

				<u>Training set</u>						
	$X_1$	$X_2$	$Y$		$X_1$	$X_2$	$Y$	$\hat{Y}$	err.	
1.	0	25	0	↗	2.	1	26	0	0	
2.	1	26	0		4.	1	27	1	1	
3.	1	22	1	↘	Error rate = $\frac{0}{2}$					
4.	1	27	1		Accuracy = 100%					
				<u>Test set</u>						
	$X_1$	$X_2$	$Y$		$X_1$	$X_2$	$Y$	$\hat{Y}$	err.	
1.	0	25	0	↗	1.	0	25	0	0	
3.	1	22	1	↘	3.	1	22	1	0	+
				Accuracy = $\frac{1}{2} = 50\% =: a_1$						



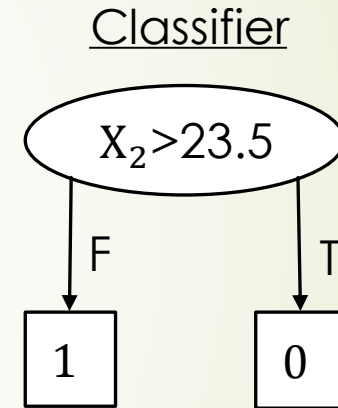
# Example of 2-fold stratified cross-validation

				<u>Training set</u>			
	$X_1$	$X_2$	$Y$	$X_1$	$X_2$	$Y$	$\hat{Y}$ err.
1.	0	25	0	1.	0	25	0 0
2.	1	26	0	3.	1	22	1 1
3.	1	22	1				
4.	1	27	1				

Error rate =  $\frac{0}{2} = 0\%$   
Accuracy = 100%

				<u>Test set</u>			
	$X_1$	$X_2$	$Y$	$X_1$	$X_2$	$Y$	$\hat{Y}$ err.
1.	0	25	0	2.	1	26	0 0
2.	1	26	0	4.	1	27	1 0 +
3.	1	22	1				
4.	1	27	1				

$$\text{Accuracy} = \frac{1}{2} = 50\% =: a_2 = a_1 = \bar{a}$$



# Deployment

- Deploy the classifier if performance is good enough, else improve the training.

- **How to deploy?**

- $N$ -tests give  $N$  classifiers. Which one to deploy?
  - Construct a final classifier using the **e**\_\_\_\_\_ data set, i.e.,  $f_{\hat{w}}$  where

$$\hat{w} \in \arg \min_{w \in \mathcal{W}} \frac{1}{|D|} \sum_{(x,y) \in D} L(y, f_w(x))$$

- **How to improve the training?**

- Consider different learning algorithms and parameters tuning.
  - Improve the pre-processing of the data.

# Controversies?

- Wouldn't the deployed classifier trained on the entire data set suffer from overfitting?
- What about deploying  $f_{\hat{w}}$  where

$$\hat{w} \in \arg \min_{w \in \mathcal{W}} \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} L(y, f_w(x)) ?$$

- Does cross validation prevent overfitting?
- Why different learning algorithms matter?

# References

- 8.5 Model Evaluation and Selection
- Optional reading:
  - Vapnik, Vladimir. "[Principles of risk minimization for learning theory](#)." *Advances in neural information processing systems*. 1992.
  - Surrogate loss functions: [https://en.wikipedia.org/wiki/Loss\\_functions\\_for\\_classification](https://en.wikipedia.org/wiki/Loss_functions_for_classification)