

# Design and Analysis of ML Experiments

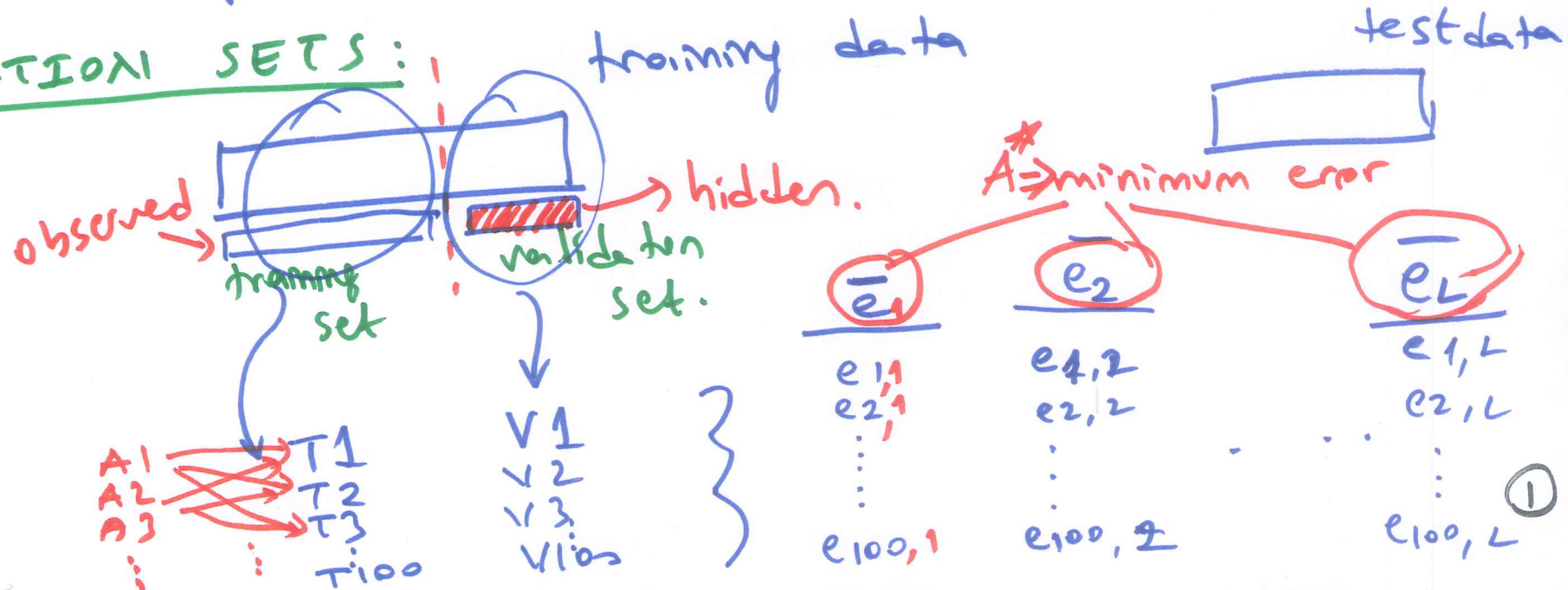
① How can we assess the expected error of a learning algorithm for a given problem?

② Given two/more algorithms, how can we say that one is better than the other(s) for a given problem?

⇒ WE CANNOT USE TRAINING SET ERRORS TO ANSWER ① and ②

training set error < ~~error~~ test set error

⇒ VALIDATION SETS:

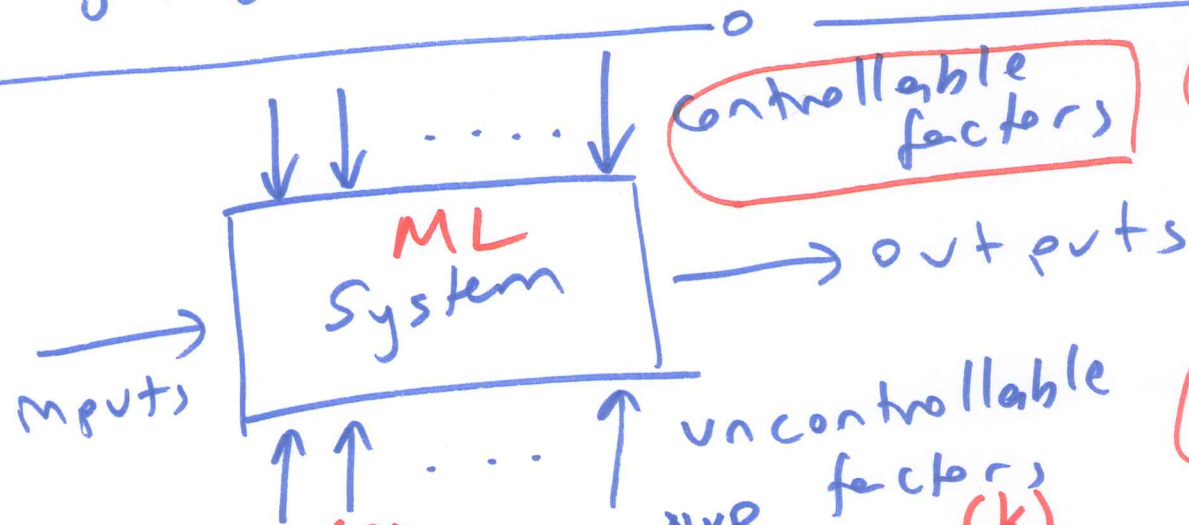


$\Rightarrow$  loss functions  $\rightarrow$  minimizing FPs  
 $\rightarrow$  minimizing FNs  
 $\rightarrow$  minimizing the cost.

$\Rightarrow$  time & space complexity.

$\Rightarrow$  interpretability.

$\Rightarrow$  easy programmability.



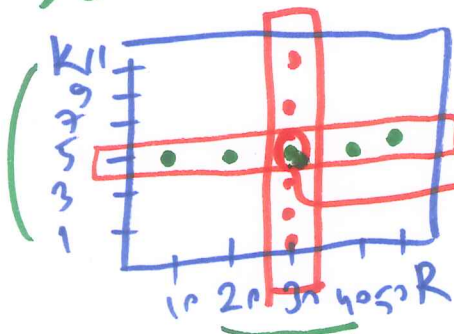
(algorithm, algorithm parameters  
input representation)

(noise in the data,  
randomness in the optimization)

$\frac{N \times D}{X} \rightarrow \text{PCA} \rightarrow \frac{N \times R}{Z} \rightarrow \text{K-NN} \rightarrow \hat{y}$

"one factor at a time"  $k$

$(R^*, k^*)$   
 "factorial design"



30 is the best given  $k=5$

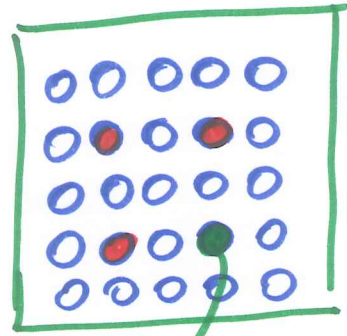
$5 + 6 - 1 = 10$  options

$5 \times 6 = 30$  options

$(R^*, k^*)$  out of  
 30 options.



## Response Surface Design.


$$\rightarrow (y_1) (e_1)$$
$$\rightarrow (y_2)(e_2)$$
$$\rightarrow (y_3)(e_3)$$

used for the model. regression algorithm

$$y_1 \rightarrow \hat{e}_1$$
 $y_2 \rightarrow e_2$ 
$$y_3 \rightarrow e_3$$

•

•

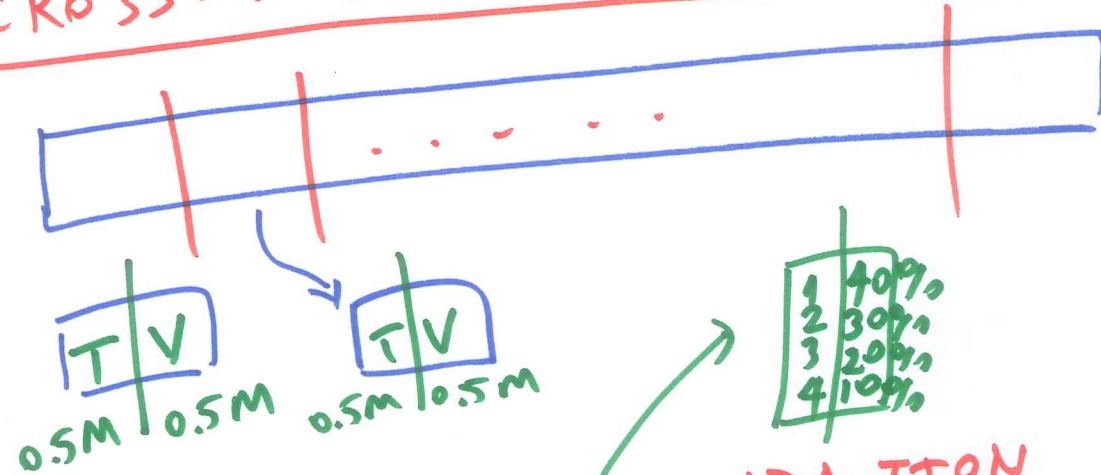
 $y_{25} \rightarrow e_{25}$ 

## GUIDELINES FOR ML EXPERIMENTS

- ## GUIDELINES FOR ML EXPERIMENT
- ① Aim of the study
    - evaluating a single algorithm and reporting its error
    - picking the best alg. on a specific problem.
    - picking the best alg. on several problems.
  - ② Selection of the response variable ] performance criteria.
  - ③ Choice of factors & their levels ]
    - algorithms
    - their parameters.
  - ④ Choice of Experimental design ]
    - do factorial design if possible
    - one factor at a time
    - response surface design.
- ③

- ⑤ Run experiments → use cloud computing if available.
- ⑥ Statistical analysis of the data results → Alg A > Alg B  
statistically significant or not?
- ⑦ Conclusions & Recommendations

## CROSS-VALIDATION & RESAMPLING



100 Million data points.

$$K=10 \Rightarrow \frac{K-2}{K-1} = \frac{8}{9}$$

## K-FOLD CROSS VALIDATION



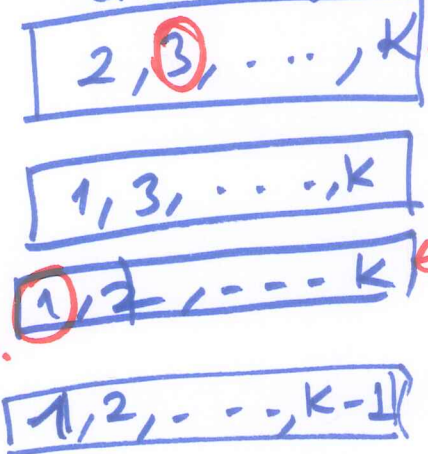
statistical correlation keeping class ratio similar between different splits.

1003  
10-fold

validation



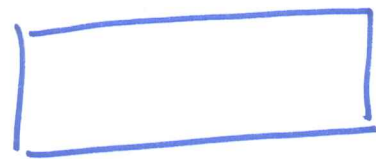
Training:





$k = N \Rightarrow$  LEAVE ONE-OUT CROSS-VALIDATION  
 40 points  $\Rightarrow$  5-fold CV  $\Rightarrow$  32 training, 8 validation  
 leave-one-out  $\Rightarrow$  39 training, 1 validation

### 5x2 CROSS VALIDATION



randomly shuffle  $\rightarrow$



$X_1^{(1)}$	$X_2^{(1)}$
-------------	-------------

$X_1^{(2)}$	$X_2^{(2)}$
-------------	-------------

$X_1^{(3)}$	$X_2^{(3)}$
-------------	-------------

$X_1^{(4)}$	$X_2^{(4)}$
-------------	-------------

$X_1^{(5)}$	$X_2^{(5)}$
-------------	-------------

