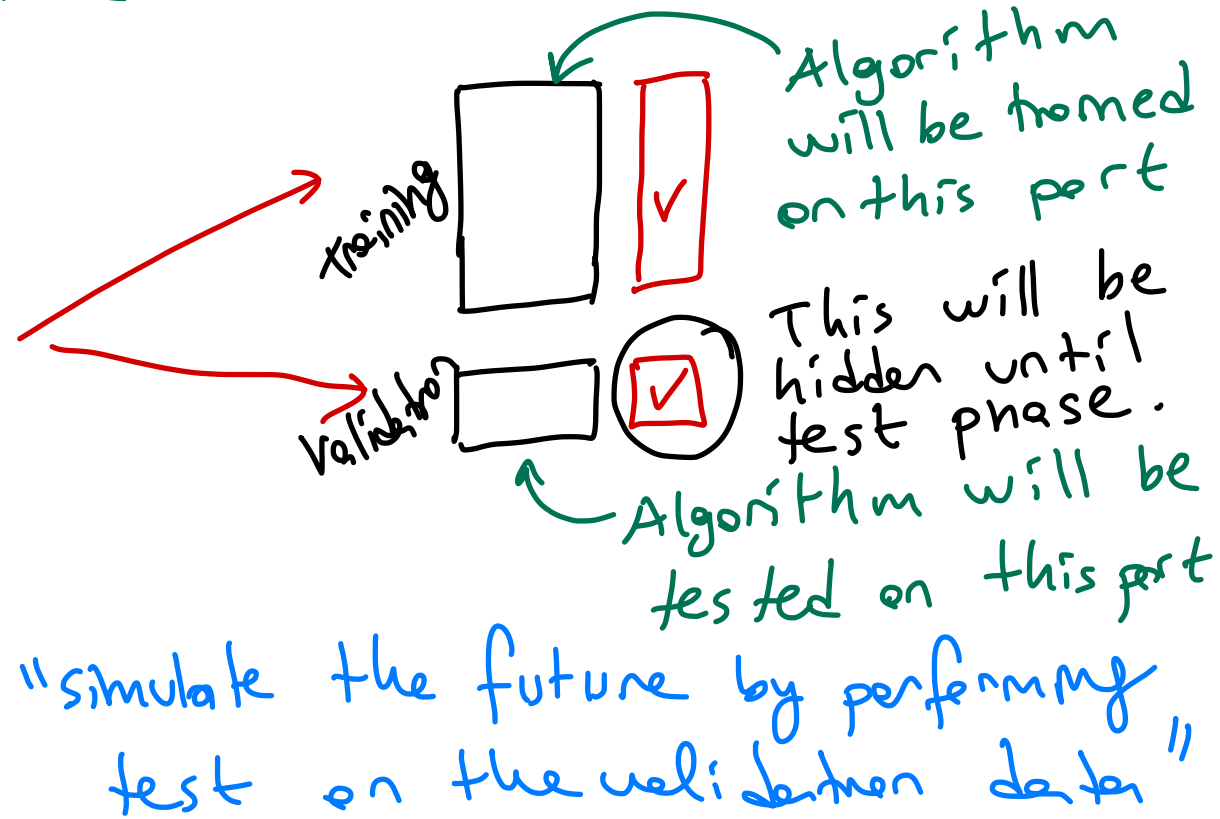
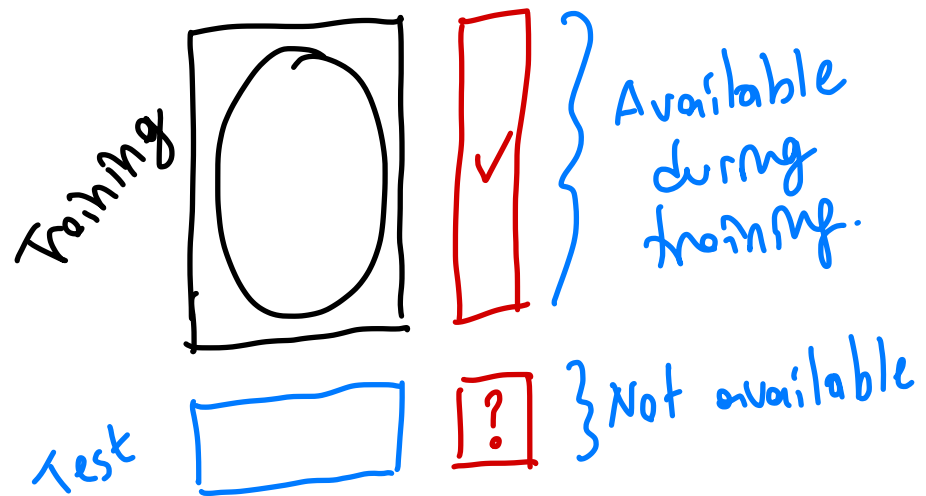


Design and Analysis of Machine Learning Experiments

- ① How can we assess the expected error of a learning algorithm for a given problem?
- ② Given two or more algorithms, how can we say that one is better than the other(s) for a given problem?

WE CAN NOT USE THE TRAINING SET TO ANSWER ① & ②.
training set error < test set error

VALIDATION SETS



A1

A2

A3

$R = \#$ of replications

R_1	T_1	V_1
R_2	T_2	V_2
R_3	T_3	V_3
\vdots	\vdots	\vdots
R_R	T_R	V_R

$e_{1,1}$ = misclassification error of A1 on R_1

$e_{3,2}$ = misclassification error of A3 on R_2

measures

- time complexity
- space complexity
- interpretability
- easy programmability
- loss functions

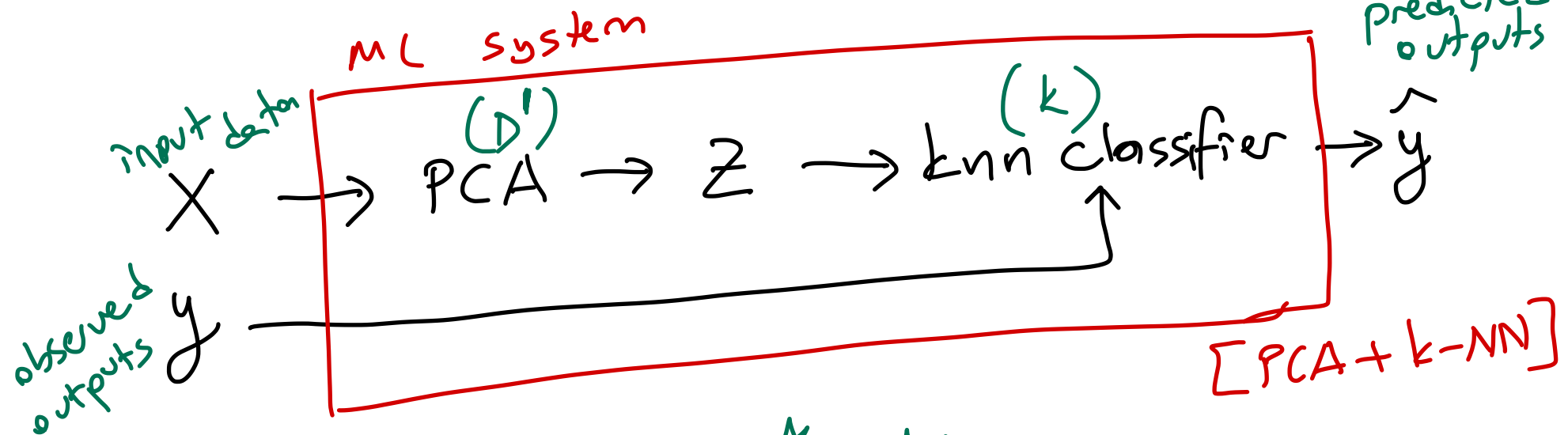
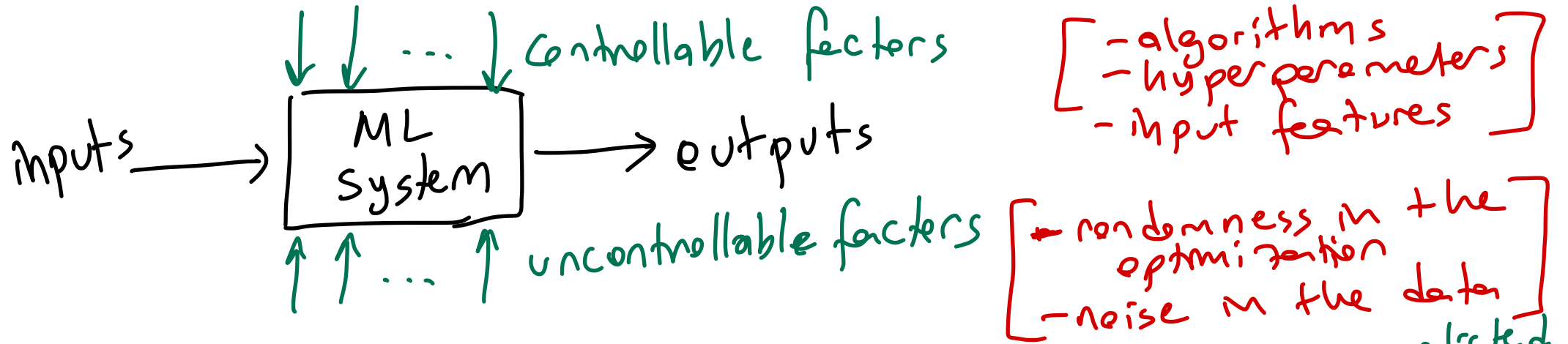
	A1	A2	A3
R_1	$e_{1,1}$	$e_{2,1}$	$e_{3,1}$
R_2	$e_{1,2}$	$e_{2,2}$	$e_{3,2}$
\vdots	\vdots	\vdots	\vdots
R_R	$e_{1,R}$	$e_{2,R}$	$e_{3,R}$
	e_1	e_2	e_3

e_1 = average performance of A1

A^* = algorithm with the minimum average error.

$A^* = A_2$

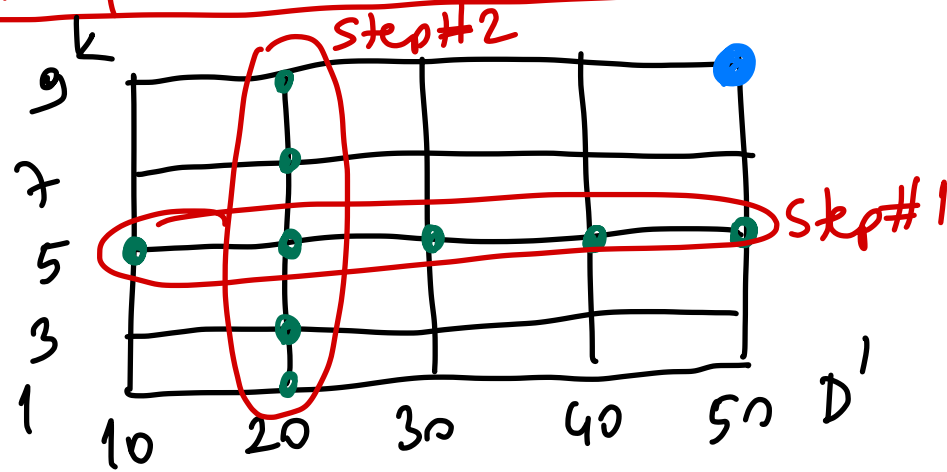
minimum



Optimization problem $\Rightarrow (D'^*, k^*)$

Exhaustive Enumeration: Try all possible combinations.
 ! Not possible due to "computational complexity"

One factor at a time:



25 possible combinations



We tried only 9 out of 25 combinations.

① Find best D' by setting k to a specified value.
Assume $k=5 \Rightarrow D'=?$

$\frac{A1}{(10,5)}$	$\frac{A2}{(20,5)}$	$\frac{A3}{(30,5)}$	$\frac{A4}{(40,5)}$	$\frac{A5}{(50,5)}$
---------------------	---------------------	---------------------	---------------------	---------------------

② Find best k by using D' from Step #1.

Assume $D'=20 \Rightarrow k=?$

$\frac{A1}{(20,1)}$	$\frac{A2}{(20,3)}$	$\frac{A3}{(20,5)}$	$\frac{A4}{(20,7)}$	$\frac{A5}{(20,9)}$
---------------------	---------------------	---------------------	---------------------	---------------------

$(D'^*, k^*) \Rightarrow (20, 7)$

Guidelines for ML Experiments

- ① Aim of the study
 - evaluate a single algorithm
 - pick the best algorithm for a specific problem
 - pick the best algorithm for a set of problems
- ② Selection of the response variable → performance criteria
- ③ Choice of factors and their levels
 - algorithms
 - hyper-parameters
 - ...
- ④ Choice of experimental design
 - ✓ exhaustive enumeration
 - factorial design
 - ✓ one factor at a time
 - response surface design
 - ...

⑤ Run experiments \rightarrow use [parallel cloud] computing if possible

⑥ Statistical analysis of the results \rightarrow Alg 1 $\stackrel{?}{>}$ Alg 2

- confidence interval
- hypothesis testing.
- ...

⑦ Conclusions & Recommendations

CROSS-VALIDATION & RESAMPLING

100 million
data points

1 million
1 million
1 million
1 million
⋮
1 million
1 million

replication #1

T_1
V_1

replication #2

T_2
V_2

⋮

replication #100

T_{100}
V_{100}

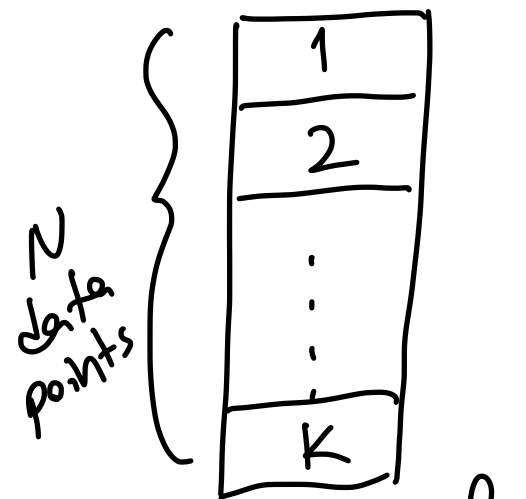
$$\frac{\Delta 1}{e_{1,1}}$$

$$\frac{A_2}{e_{2,1}}$$

$$\frac{e_{1,2}}{e_{2,2}}$$

$$\frac{e_{1,100}}{e_1} \stackrel{?}{>} \frac{e_{2,100}}{e_2}$$

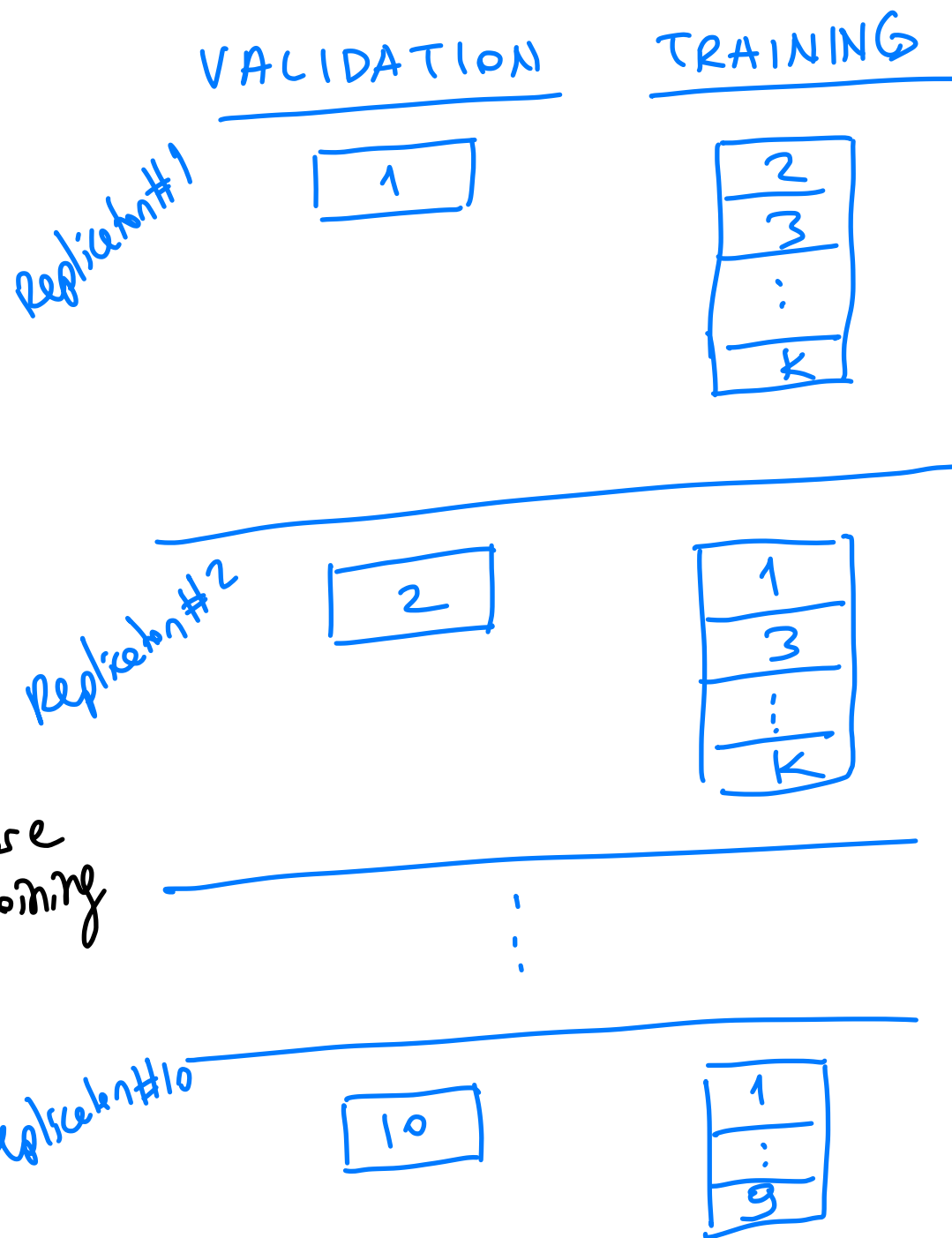
K-FOLD CROSS-VALIDATION



almost equal
"K" partitions

- * No overlap in validation!
- * $(K-2)$ out of $(K-1)$ blocks are common in training

Overlap between training data is quite large

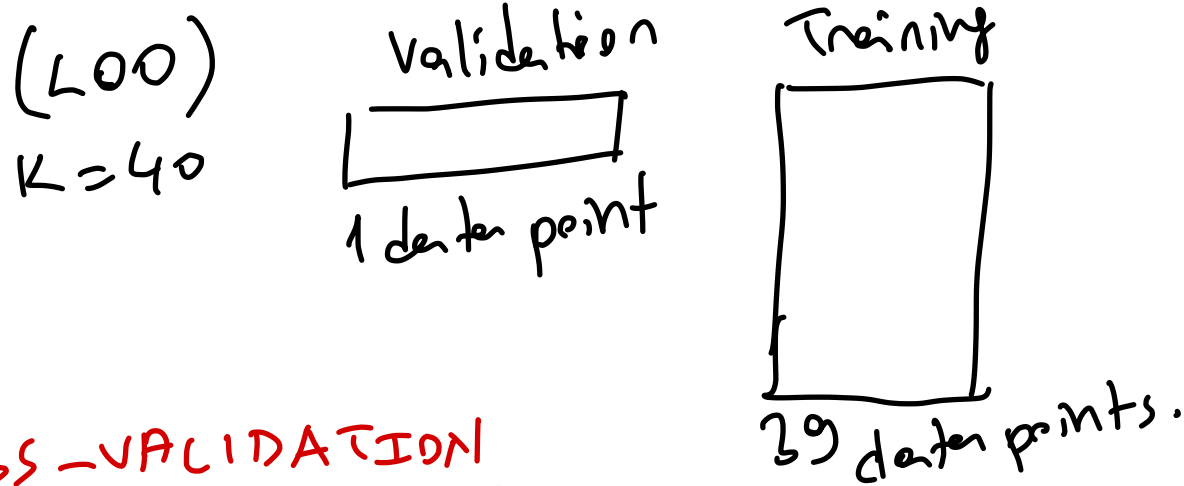
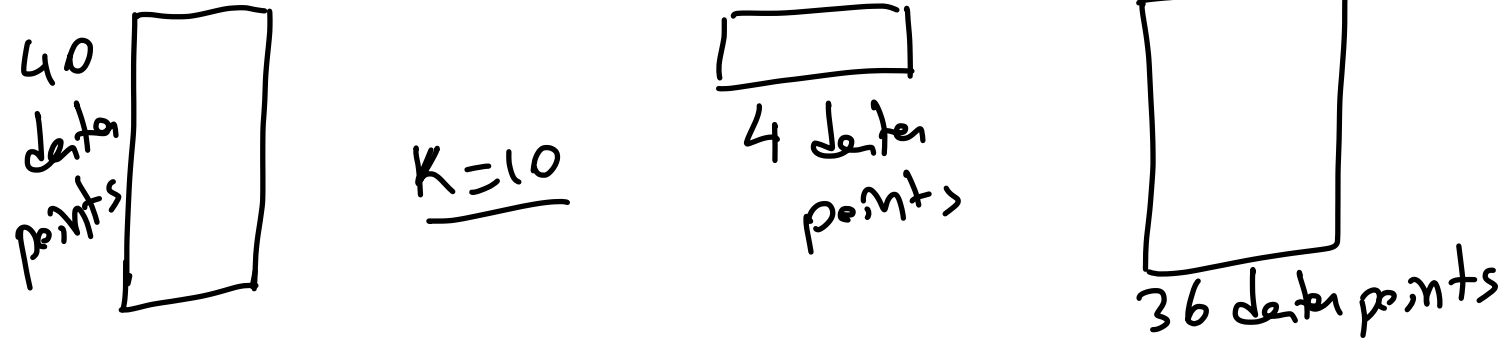


Step 1: Shuffle data points.

Step 2: Divide the data set into K (almost) equal sized partitions

Step 3: Use each partition as the validation data set in one replication

LEAVE-ONE-OUT CROSS-VALIDATION ((K-FOLD cross validation where $K=N$))



5x2 CROSS-VALIDATION

