

Introduction To Machine Learning

- Introduction and Machine Learning Principles
- Examples of Classifiers
- How to Perform Experiments
- Metrics
- Suggested Readings



Machine learning studies computer algorithms for learning to do stuff.
(Schapire)

A computer program is set to learn from an experience E with respect to some task T and some performance measure P if its performance on T as measured by P improves with experience E .
(Mitchell)

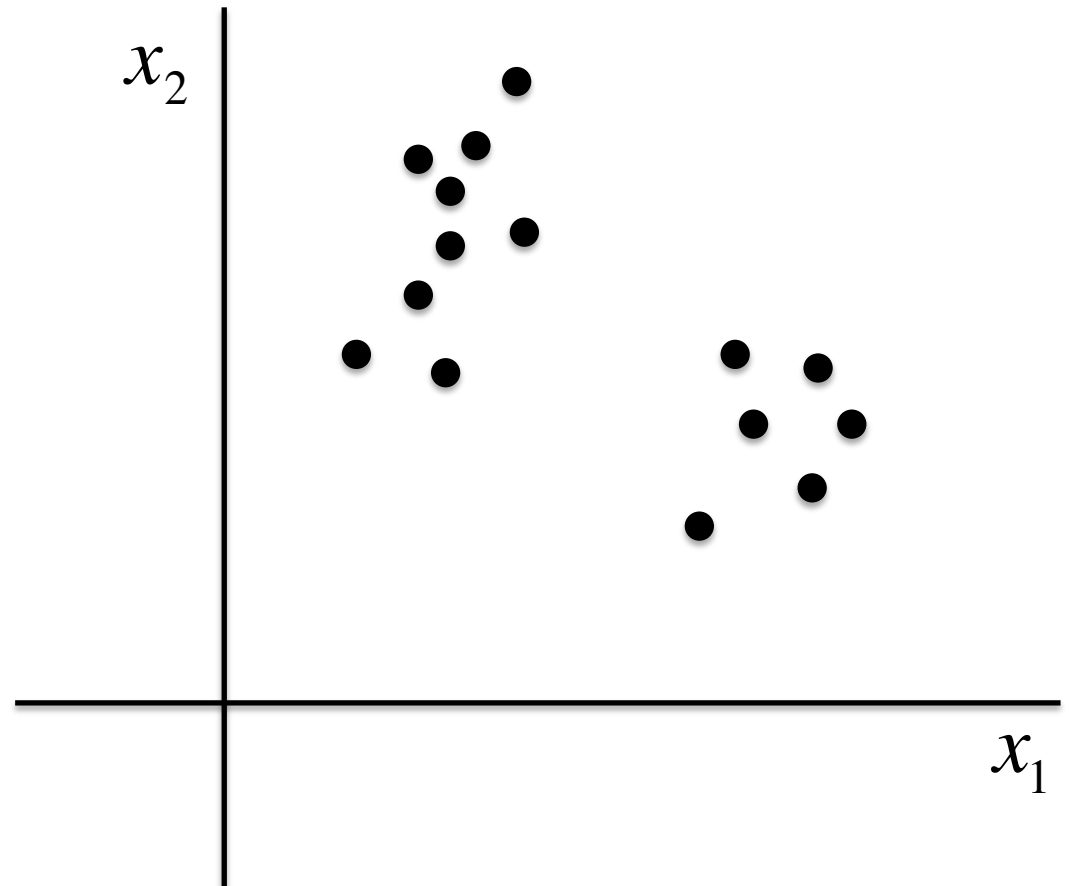
$$\bar{x} = (x_1, \dots, x_m)^T$$

Sample

$$\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N\}$$

Dataset

$$\bar{x} = (x_1, x_2)^T$$



$$\bar{x} = (x_1, \dots, x_m)^T$$

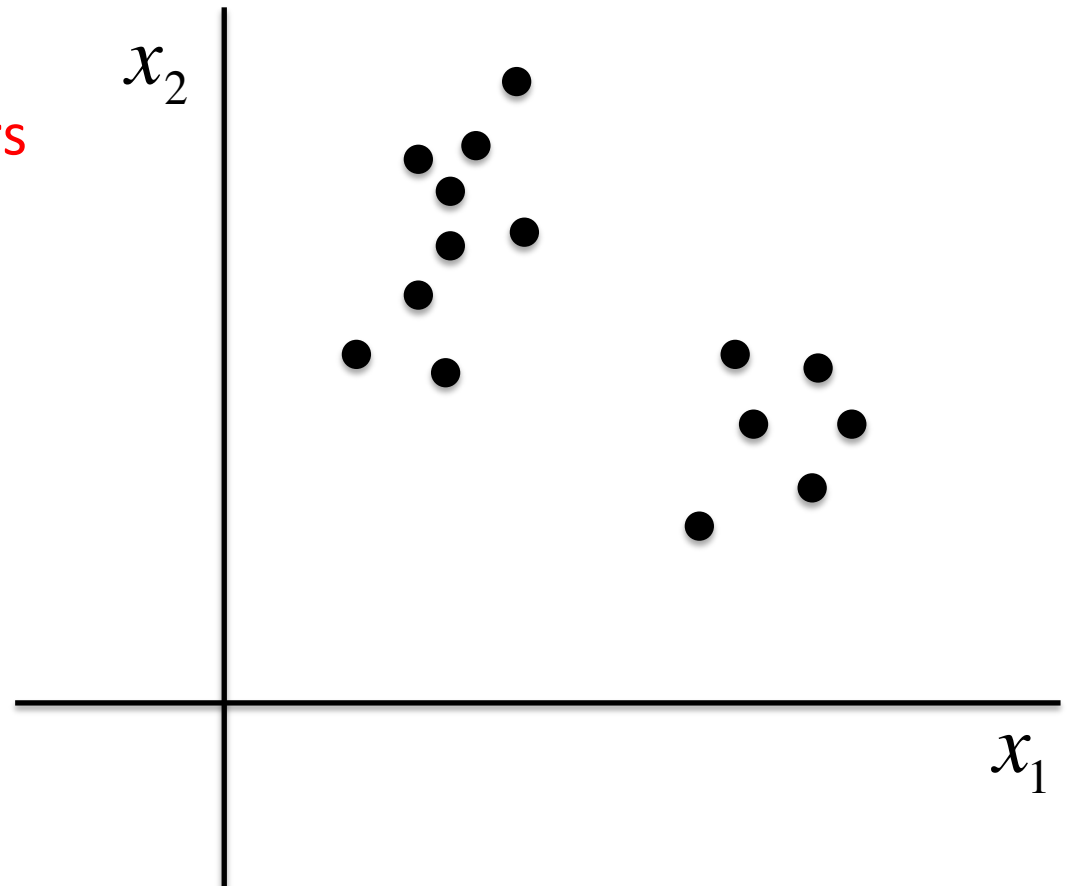
Sample

Features or descriptors

$$\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N\}$$

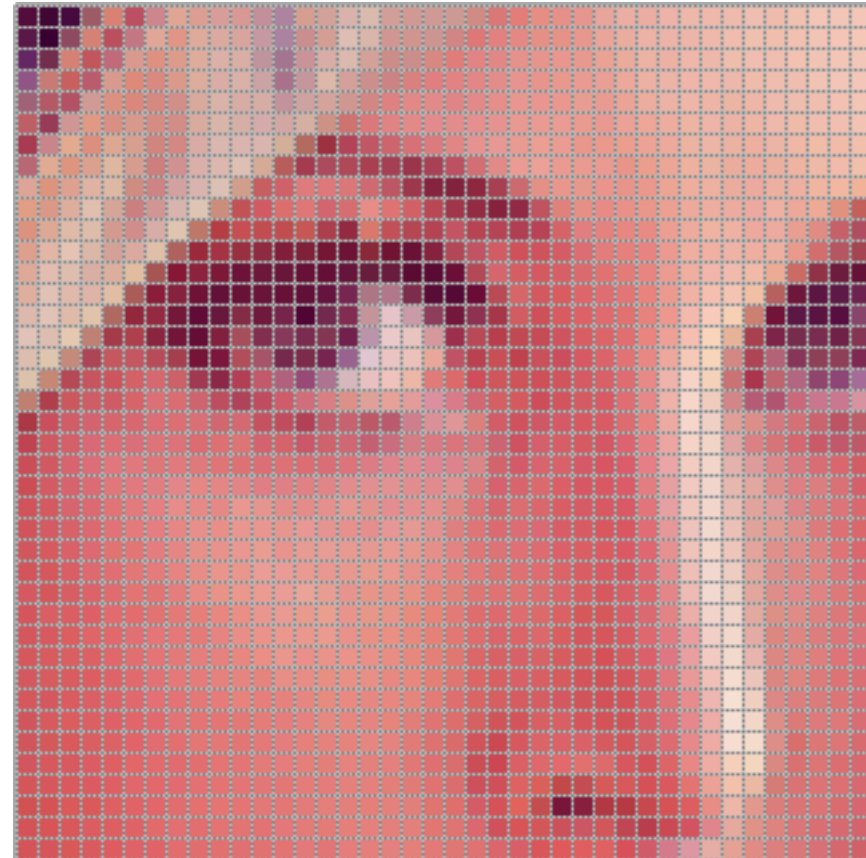
Dataset

$$\bar{x} = (x_1, x_2)^T$$



$$\bar{x} = (x_1, \dots, x_m)^T$$

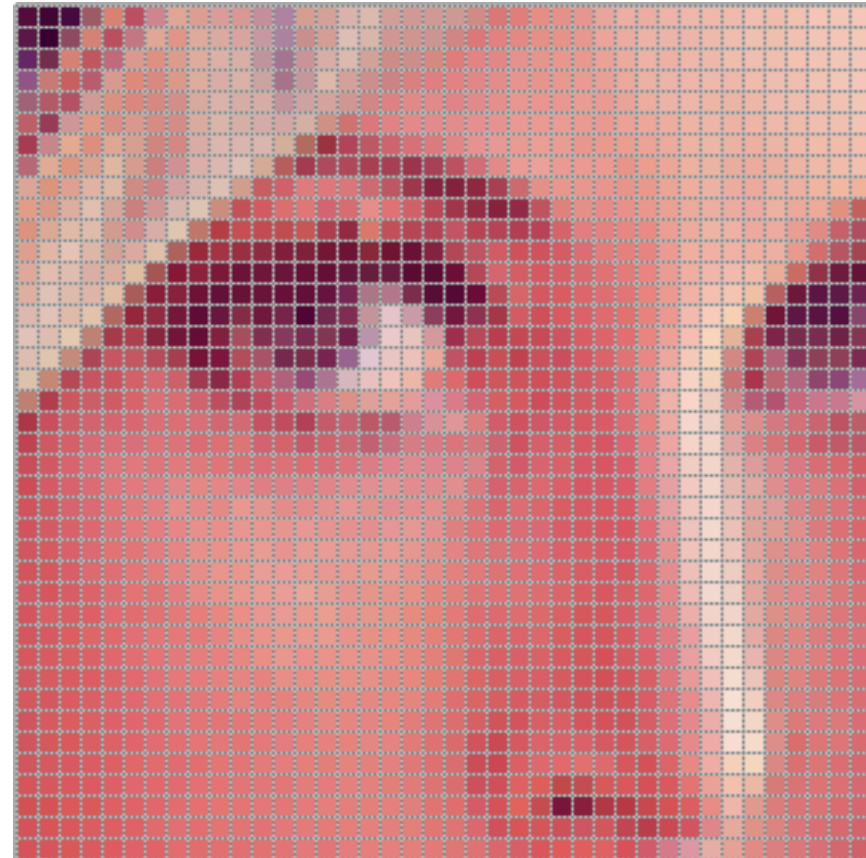
Sample



$$\bar{x} = (x_1, \dots, x_m)^T$$

Sample

1 Mpx $\Rightarrow 10^6$ features

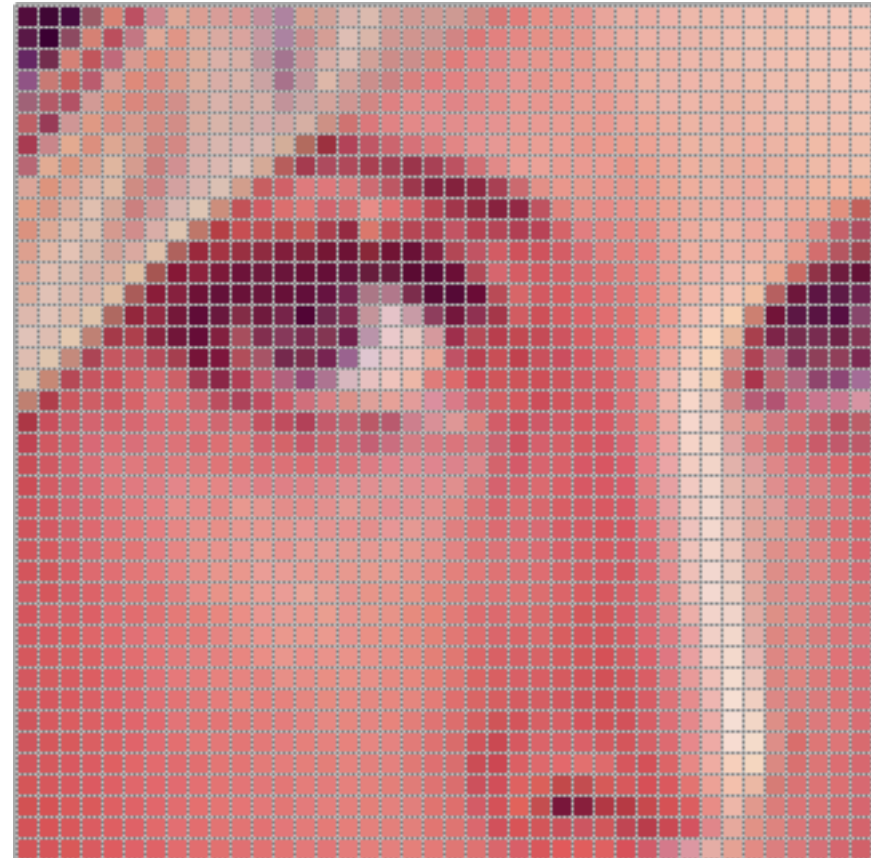


$$\bar{x} = (x_1, \dots, x_m)^T$$

Sample

1 Mpx $\Rightarrow 10^6$ features

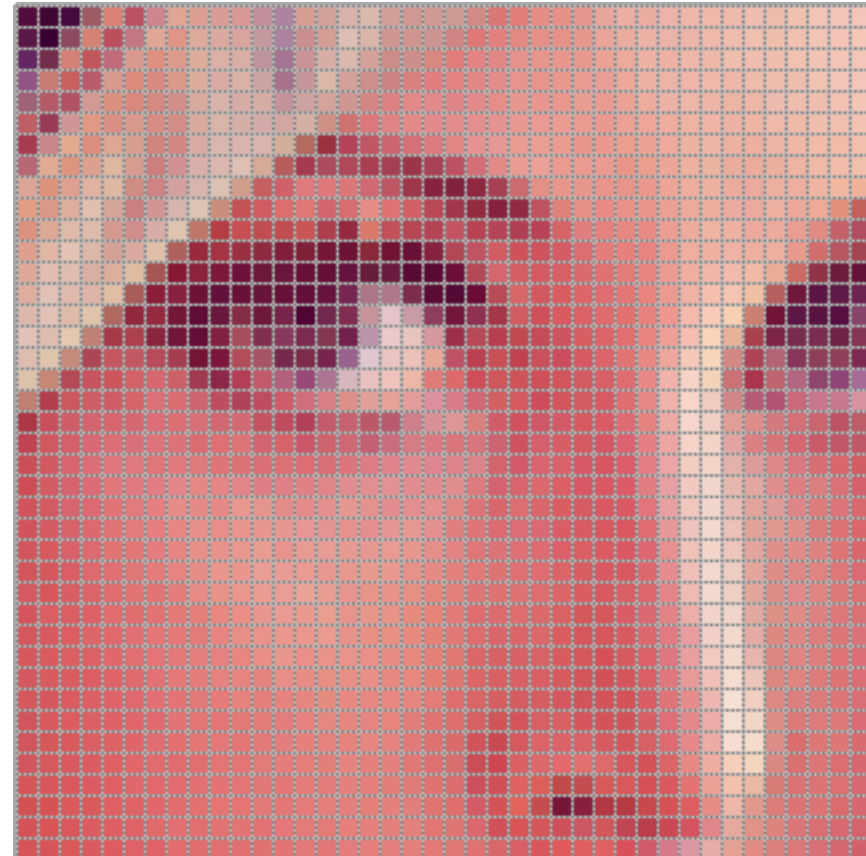
- Computationally expensive
- Curse of dimensionality



$$\bar{x} = (x_1, \dots, x_m)^T$$

Sample

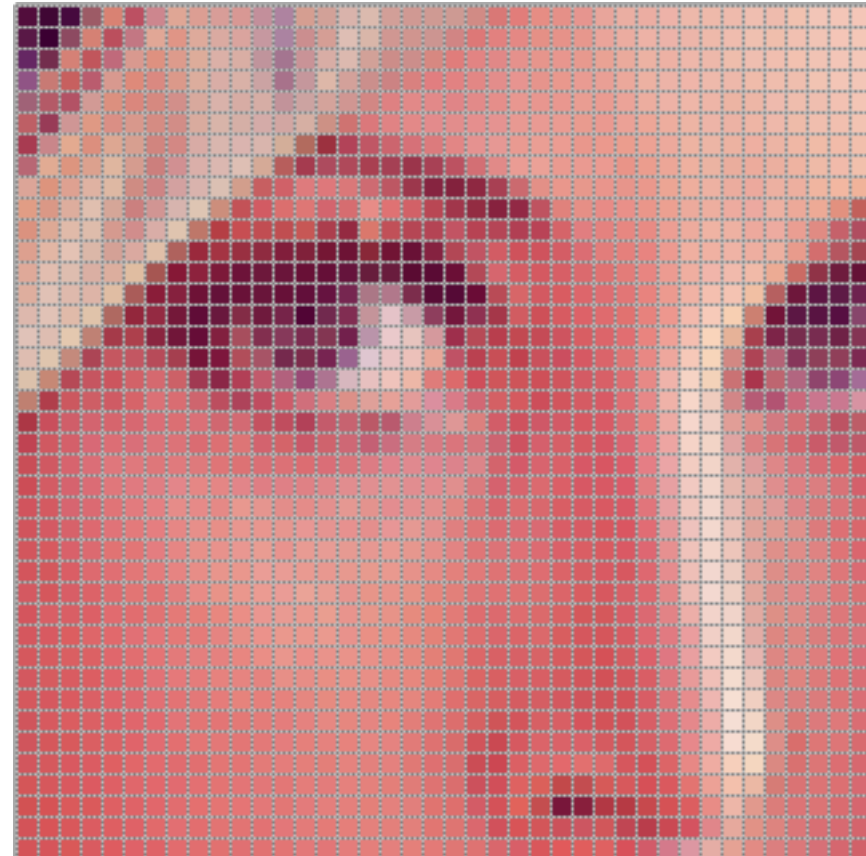
- Mean, variance, entropy,...
- SIFT
- Harris Corner
-



$$\bar{x} = (x_1, \dots, x_m)^T$$

Sample

Feature extraction is essentially a dimensionality reduction problem with the goal of finding meaningful projections of the original data vectors

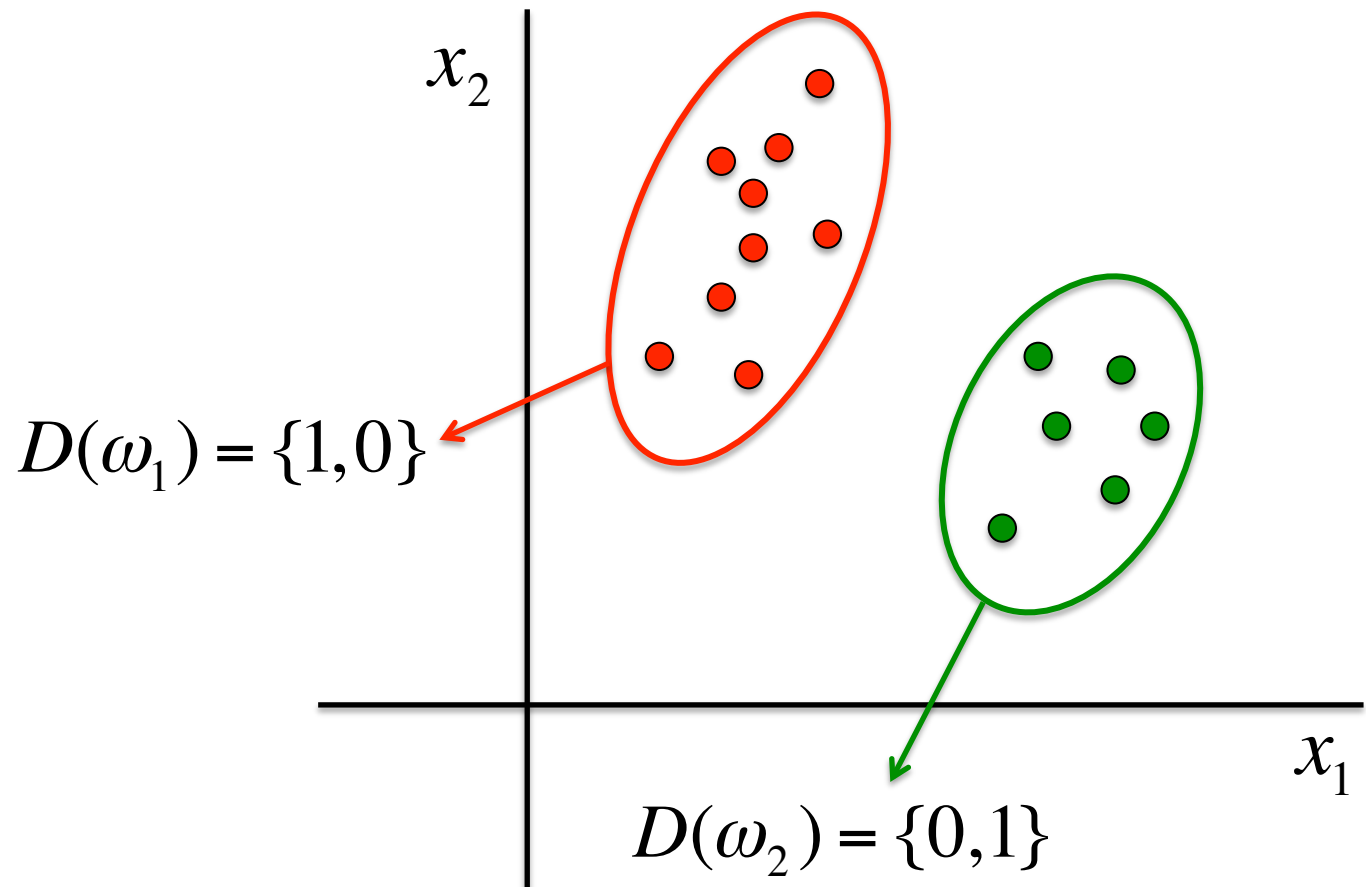


$$\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$$

Class

$$D(\omega_i) = \{0, 0, \dots, 1, \dots, 0\}$$

Label



CLASSIFICATION PROBLEM

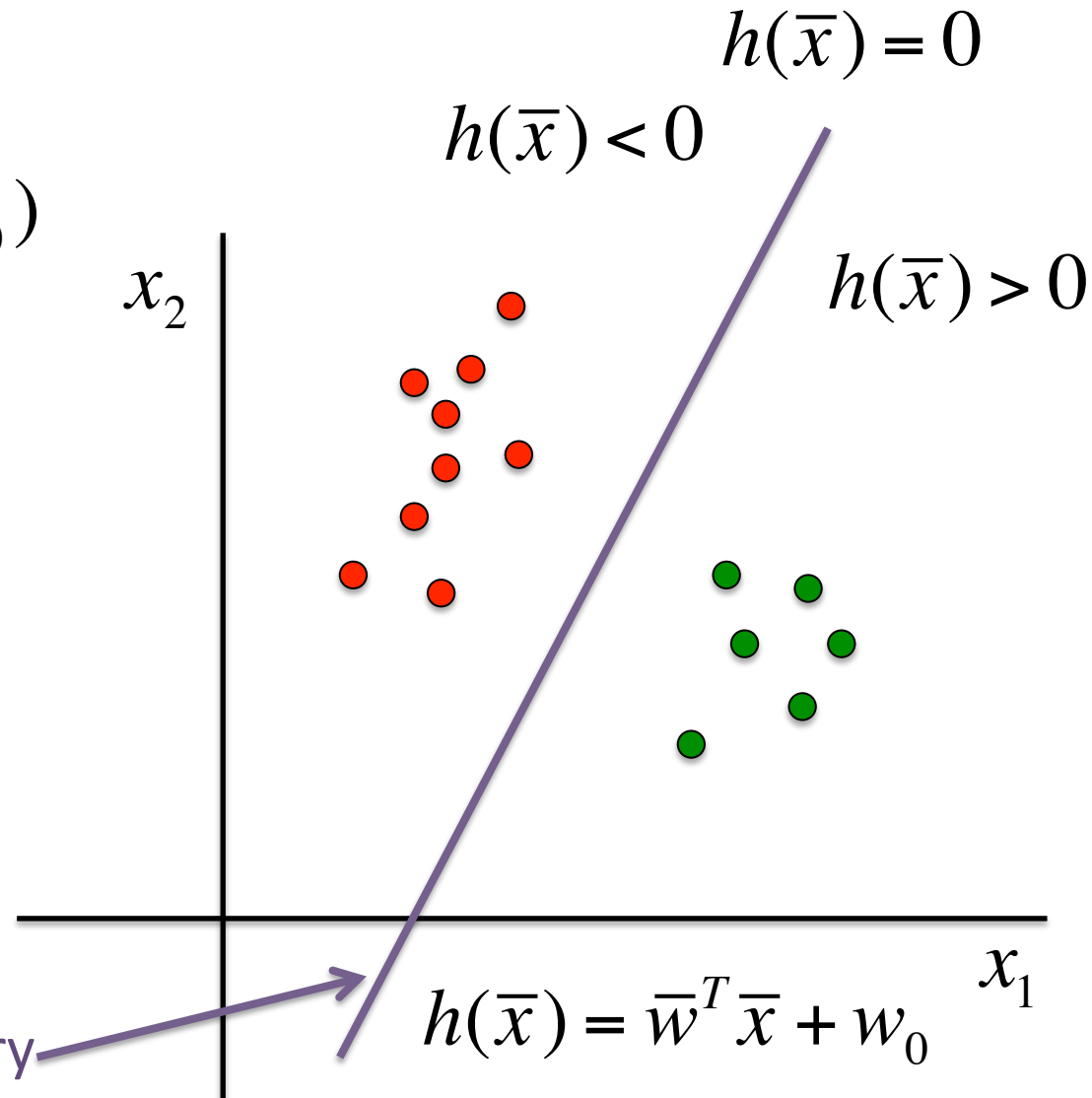
$$h(\bar{x}) = f(\bar{w}^T \bar{x} + w_0)$$

Classifier

$$f(\cdot)$$

Activation function

Decision Boundary

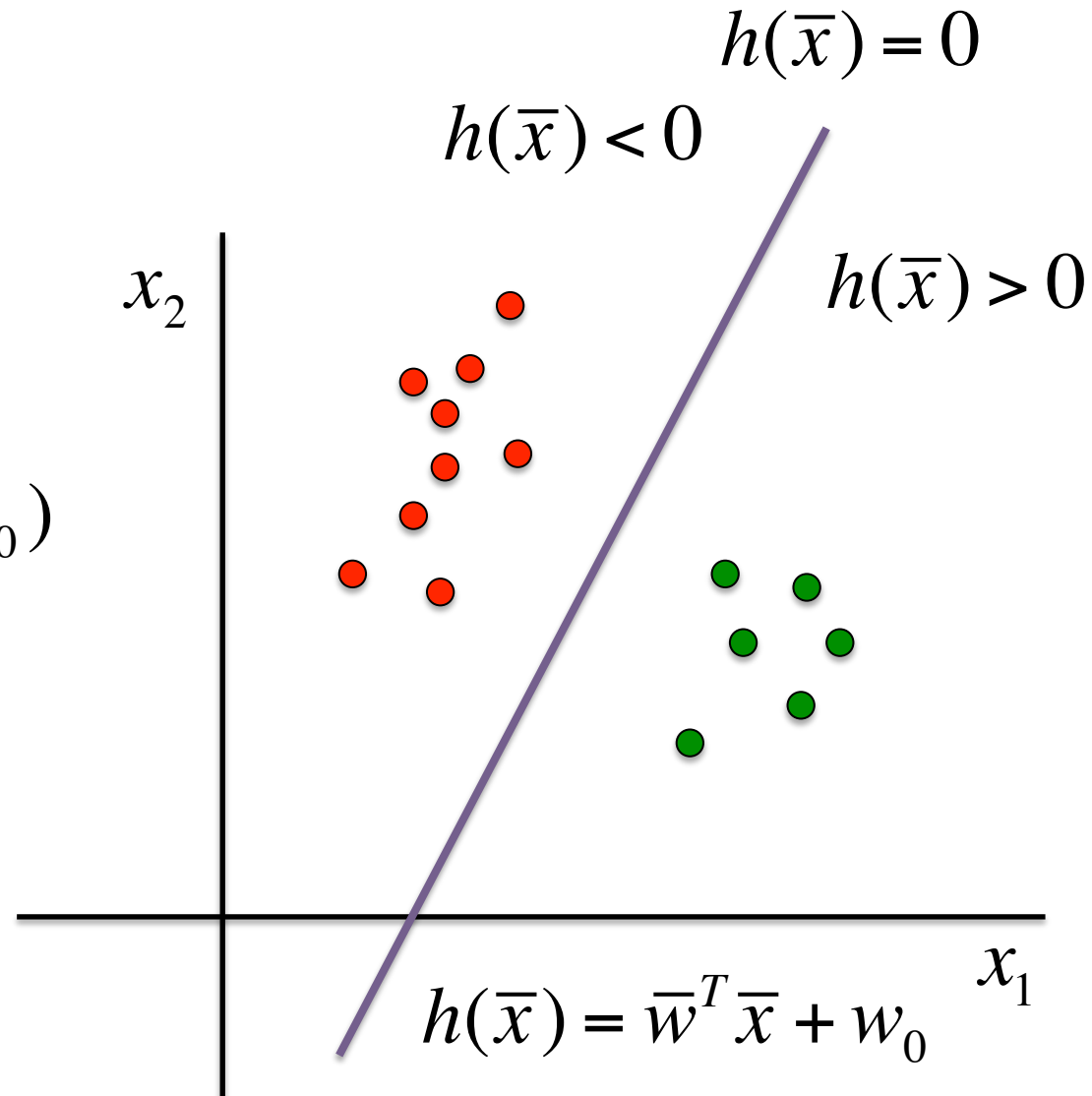


$$h(\bar{x}) = f(\bar{w}^T \bar{x} + w_0)$$

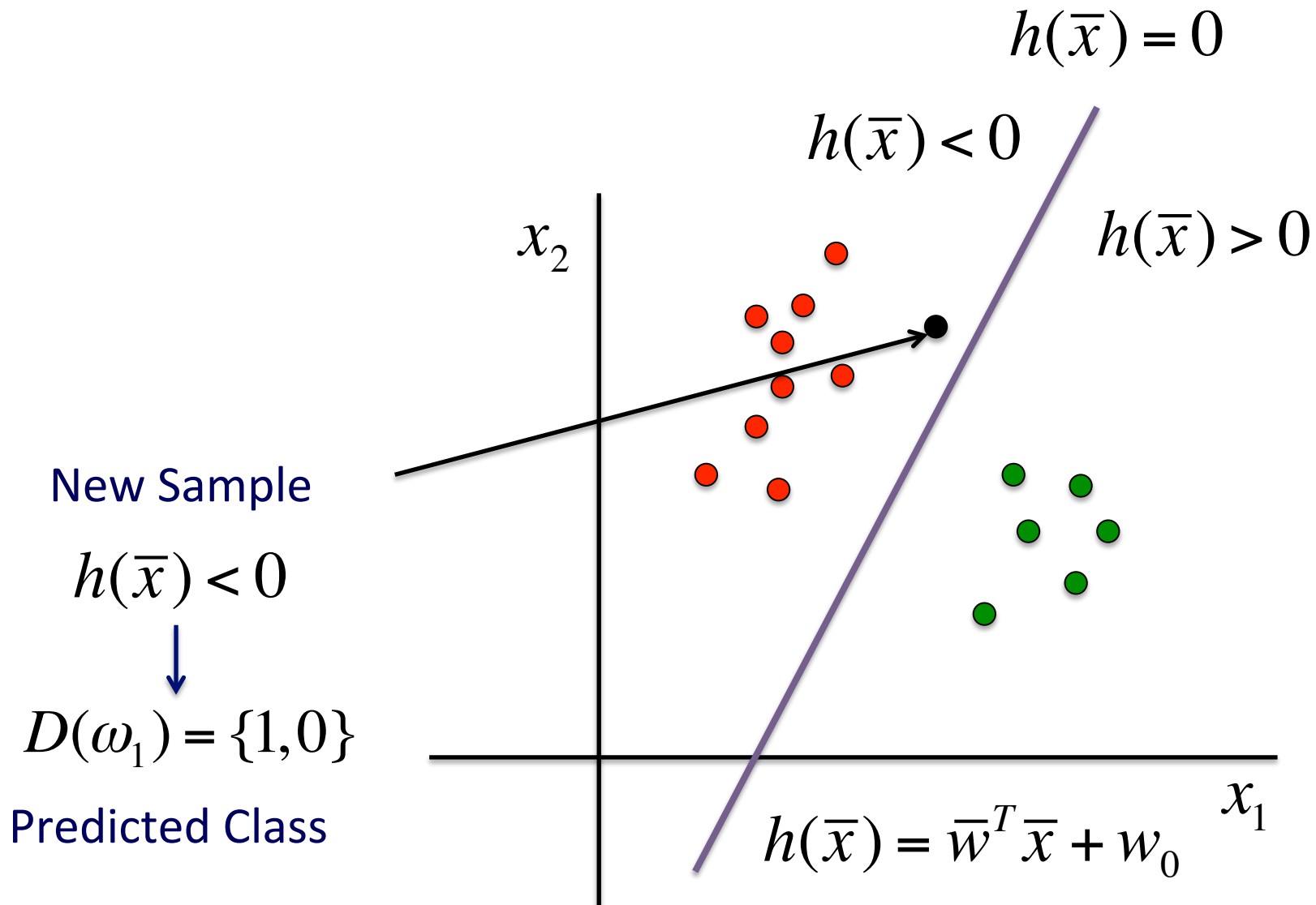


$$L(\omega(\bar{x}), h(\bar{x}))$$

Loss Function

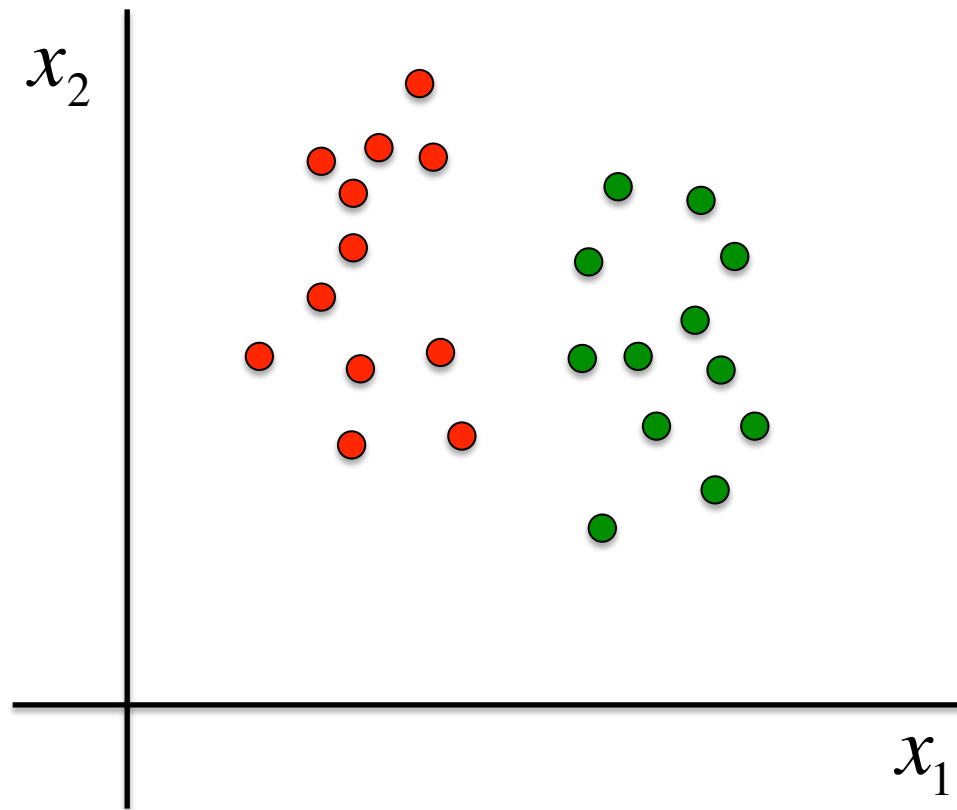


CLASSIFIER TESTING



- k Nearest Neighbors (kNN)
- Neural Network (NN)
- Support Vector Machine

k Nearest Neighbors

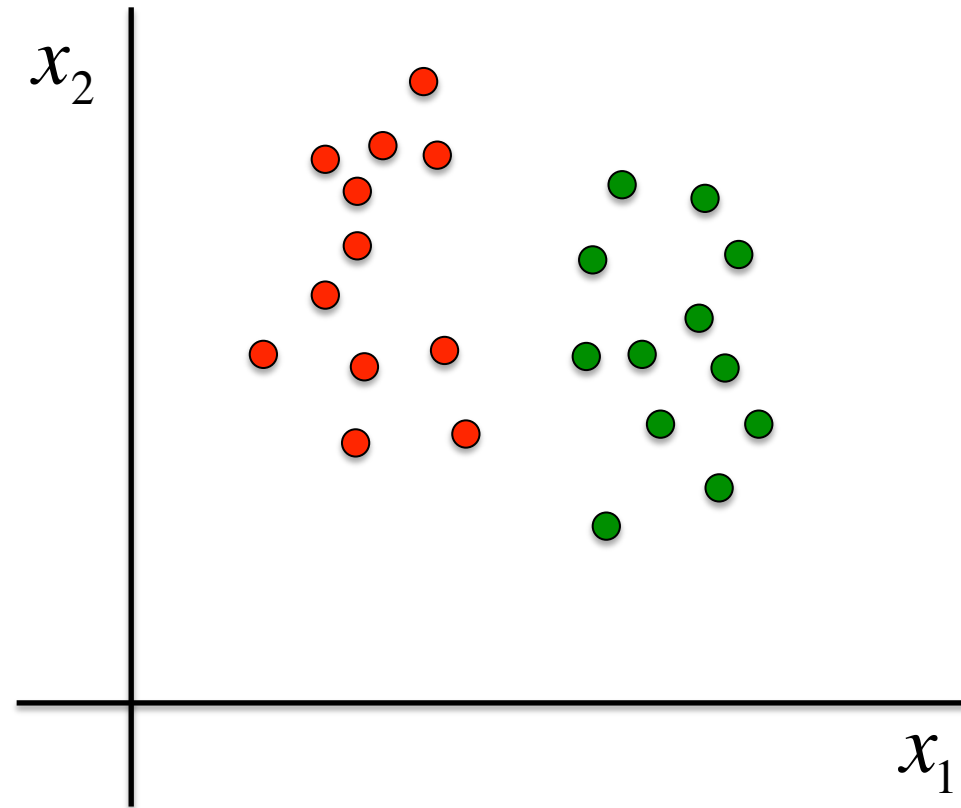


k Nearest Neighbors

No training set

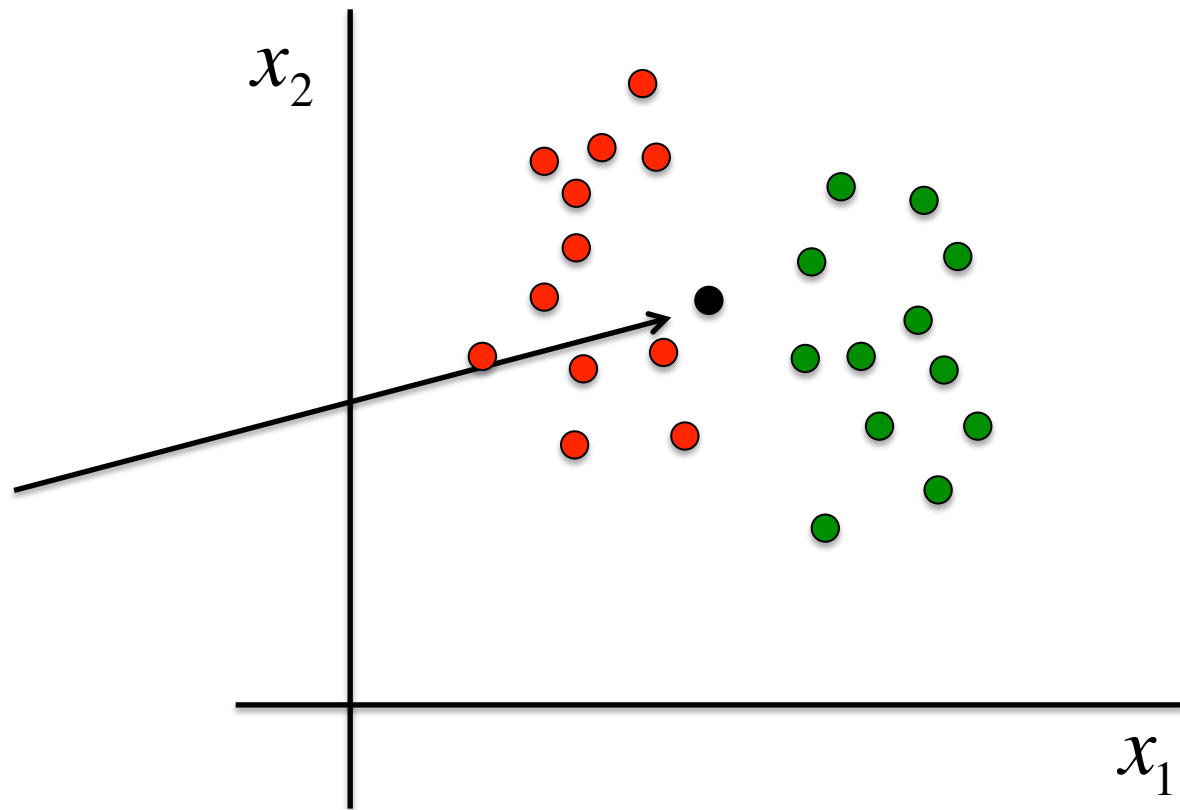


Reference Set



k Nearest Neighbors

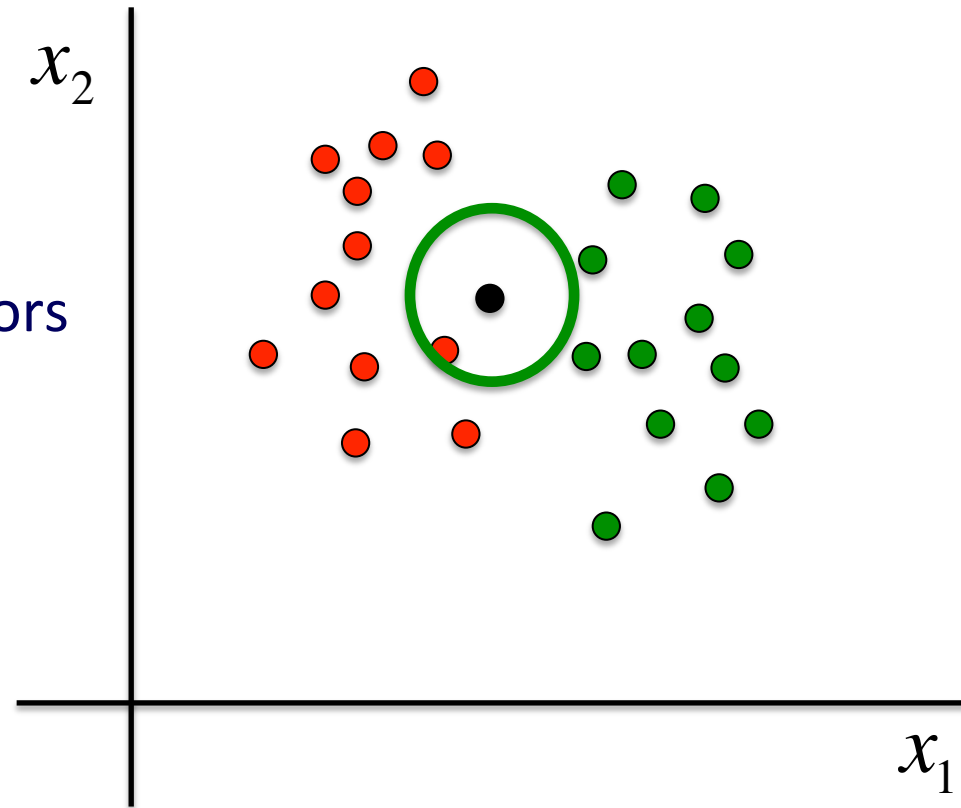
New Sample



k Nearest Neighbors

Algorithm steps:

- Consider the fixed value k
e.g. $K=1$
- Find the K nearest neighbors
- Set the class according to
a majority voting



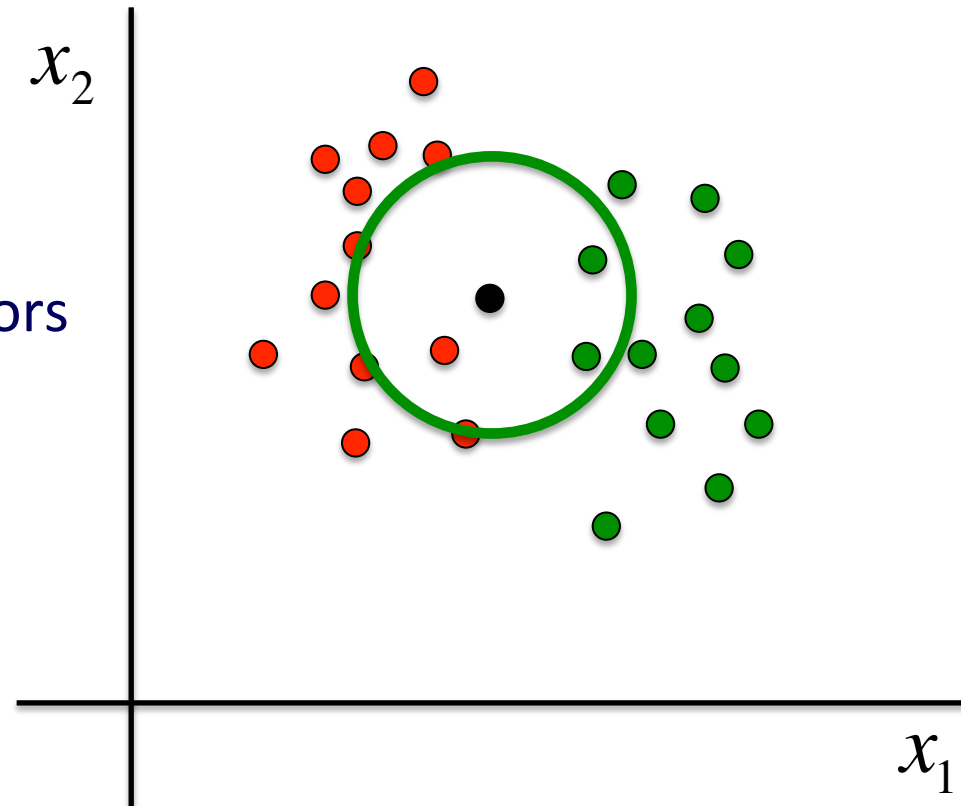
k Nearest Neighbors

- Is one of the simplest machine learning algorithm;
- There is no training stage;
- The classification function is approximated locally;
- Several variants of the original algorithm have been proposed.

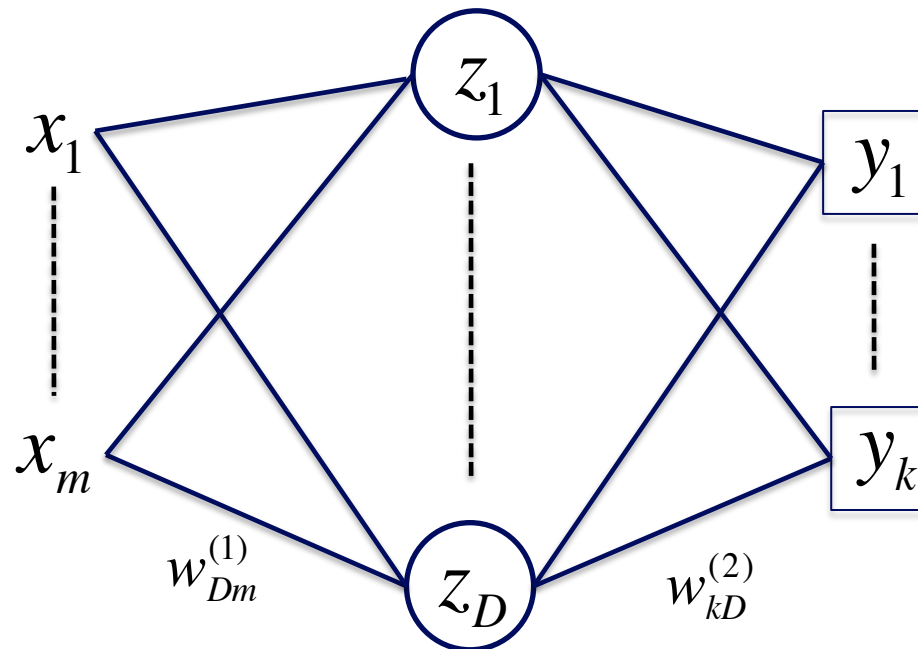
k Nearest Neighbors

Algorithm steps:

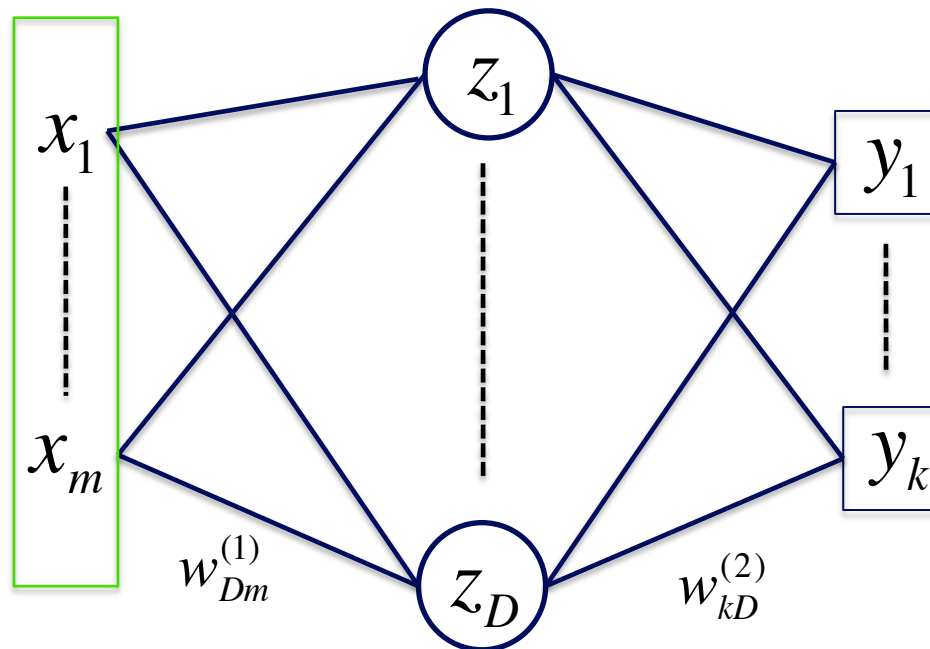
- Consider the fixed value k
e.g. $K=3$
- Find the K nearest neighbors
- Set the class according to a majority voting



Inspired by the central nervous system, a Neural Network is a set of interconnected “neurons” which compute value from inputs

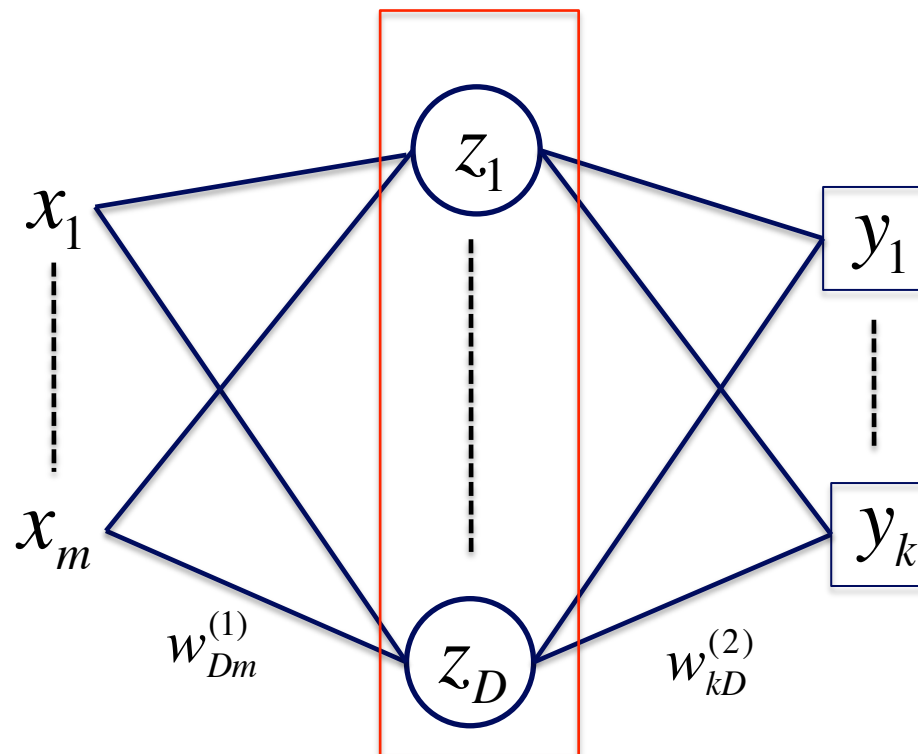


Inspired by the central nervous system, a Neural Network is a set of interconnected “neurons” which compute value from inputs



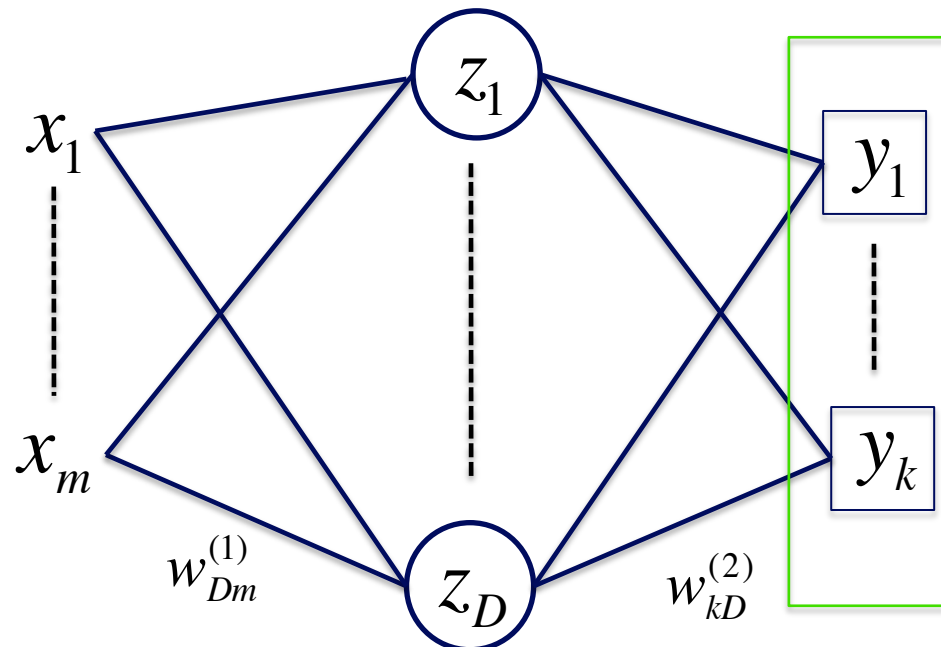
Input Layer

Inspired by the central nervous system, a Neural Network is a set of interconnected “neurons” which compute value from inputs



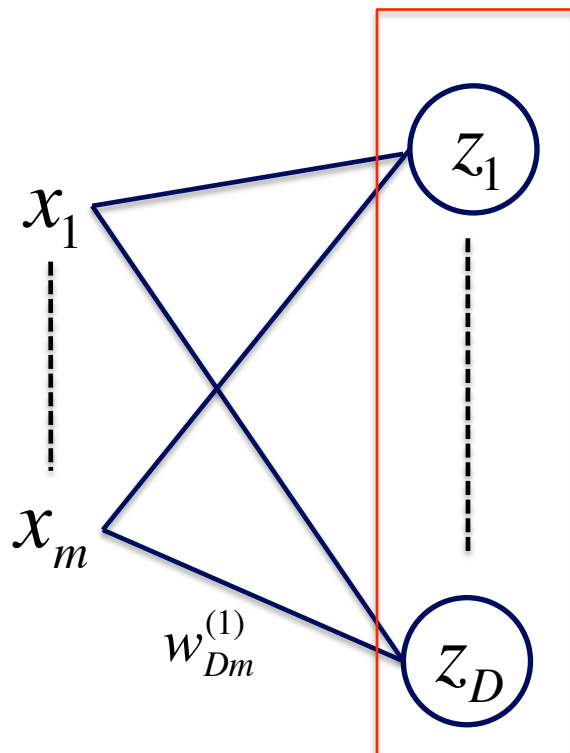
Hidden Layer

Inspired by the central nervous system, a Neural Network is a set of interconnected “neurons” which compute value from inputs



Output Layer

The value of each hidden unit is computed by an activation function according to the value of weights of the first layer and the inputs.



Hidden Layer

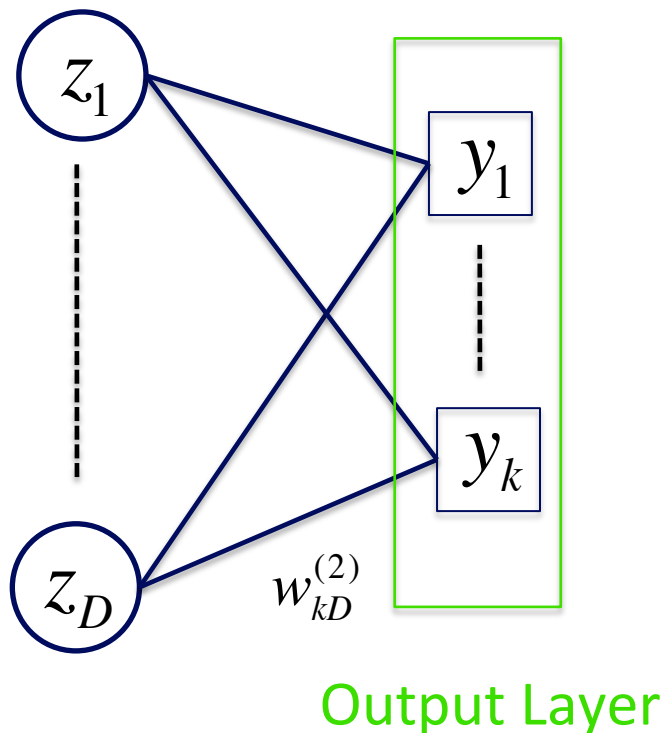
$$z_j = h(a_j)$$

Activation function

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i$$

Input Linear
Combination

The output value is computed by an activation function according to the values of the weights of the second layer and the hidden units.



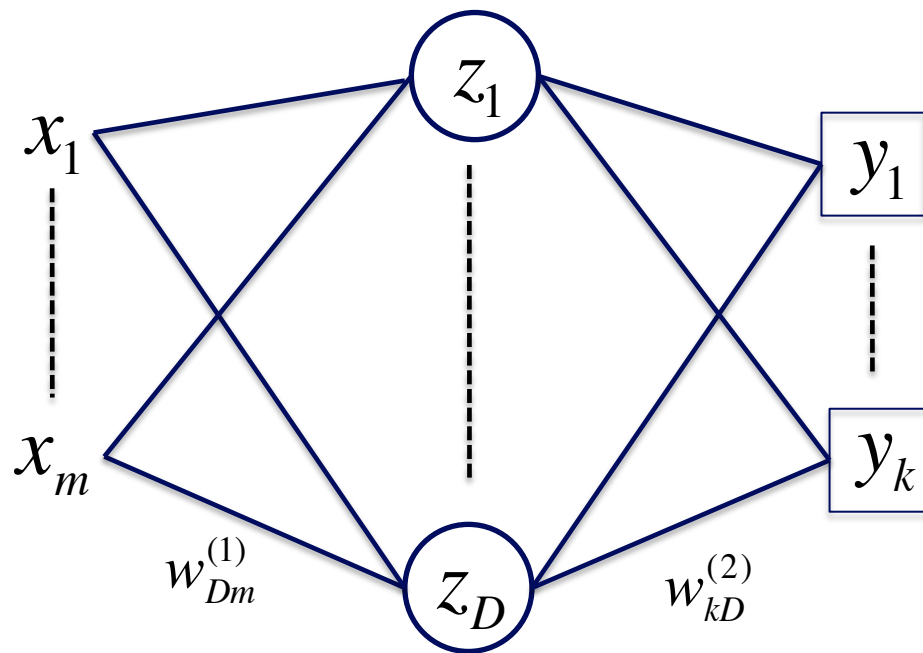
$$y_k = \sigma(a_k)$$

Activation function

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j$$

Output unit
activations

The training of the network consists in estimating the layers' weights



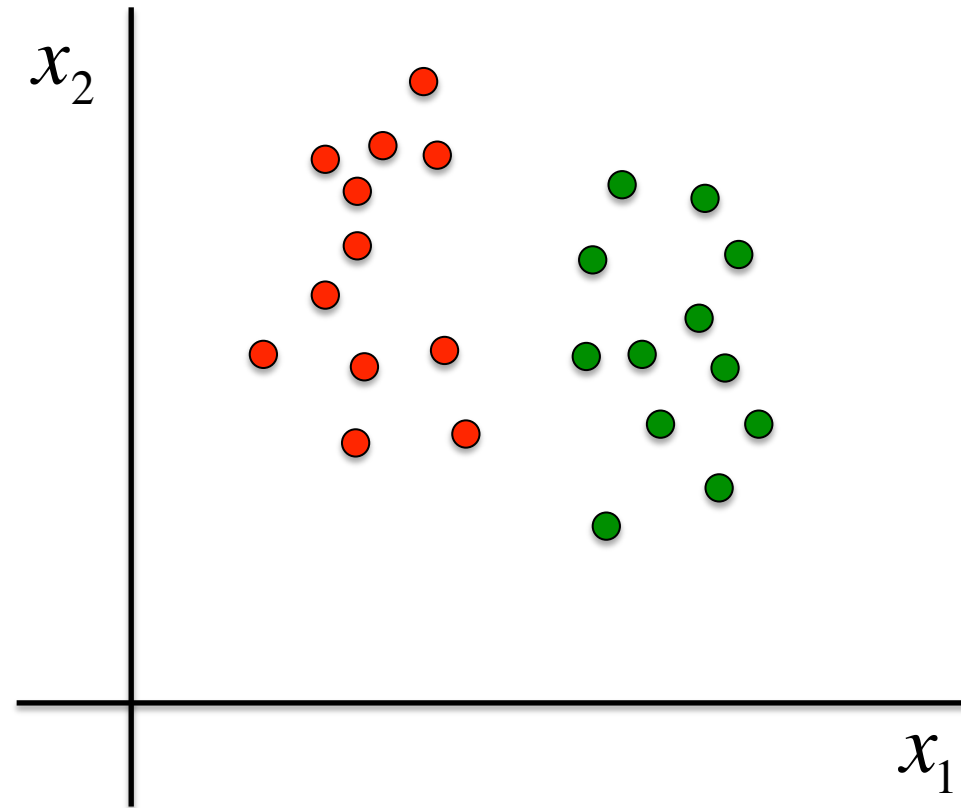
$$E(w) = \frac{1}{2} \sum_{n=1}^N \left\| y(\bar{x}, w) - D(\omega) \right\|^2$$

- Activation function generally are sigmoid function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

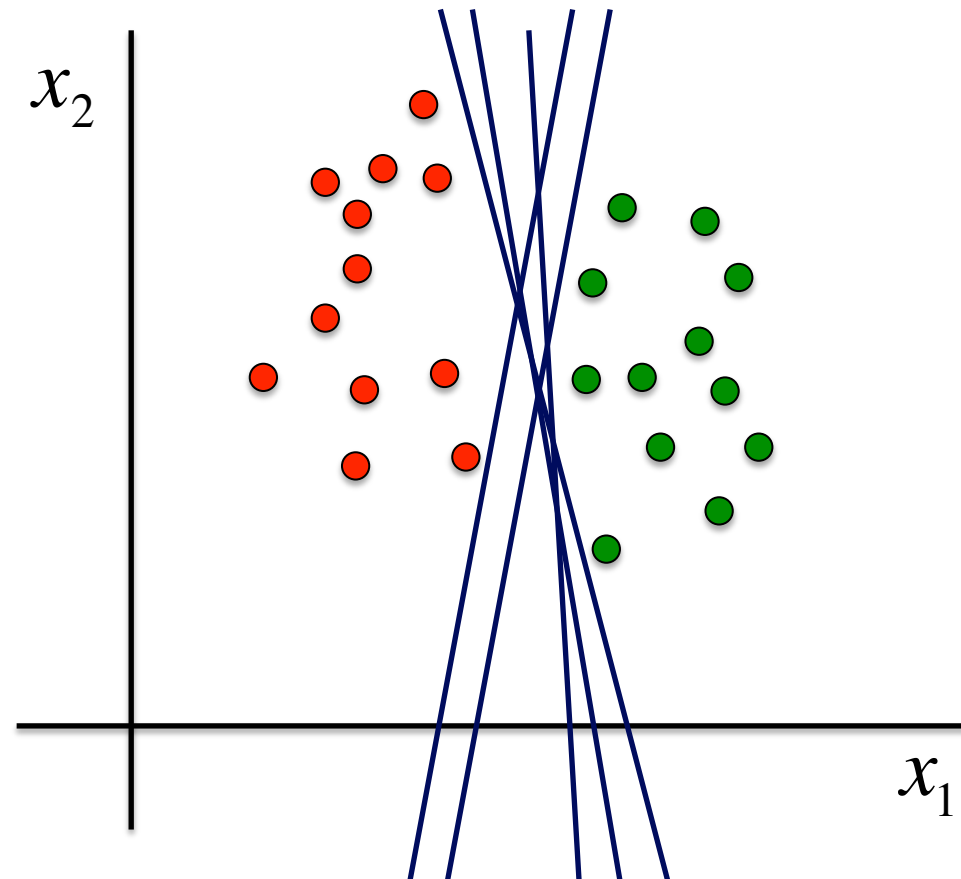
- The number of layers and the number of the hidden units are hyper-parameters of the model

In the support vector machines the decision boundary is chosen to be the one for which the margin is maximized.



In the support vector machines the decision boundary is chosen to be the one for which the margin is maximized.

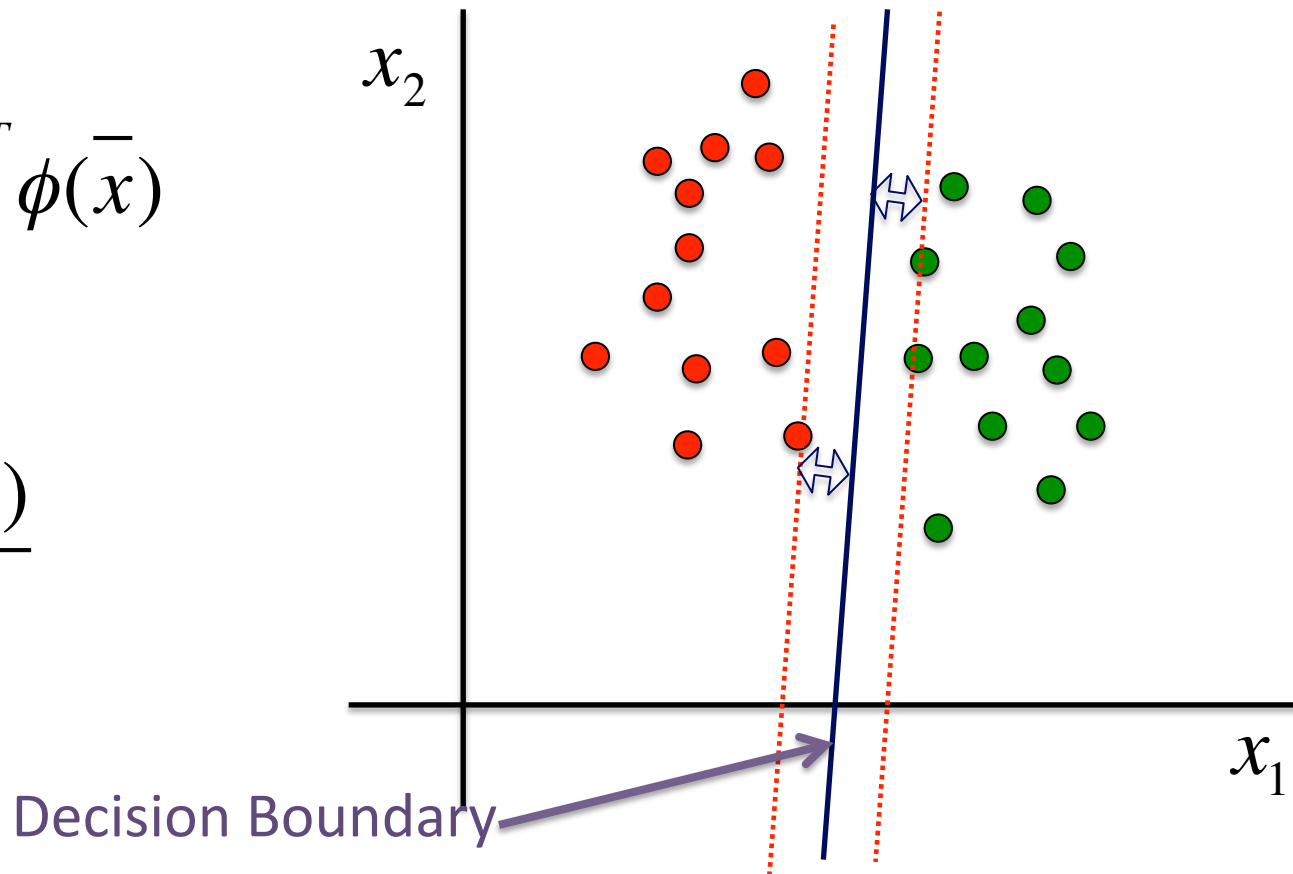
$$y(\bar{x}_n) = \bar{w}^T \phi(\bar{x})$$



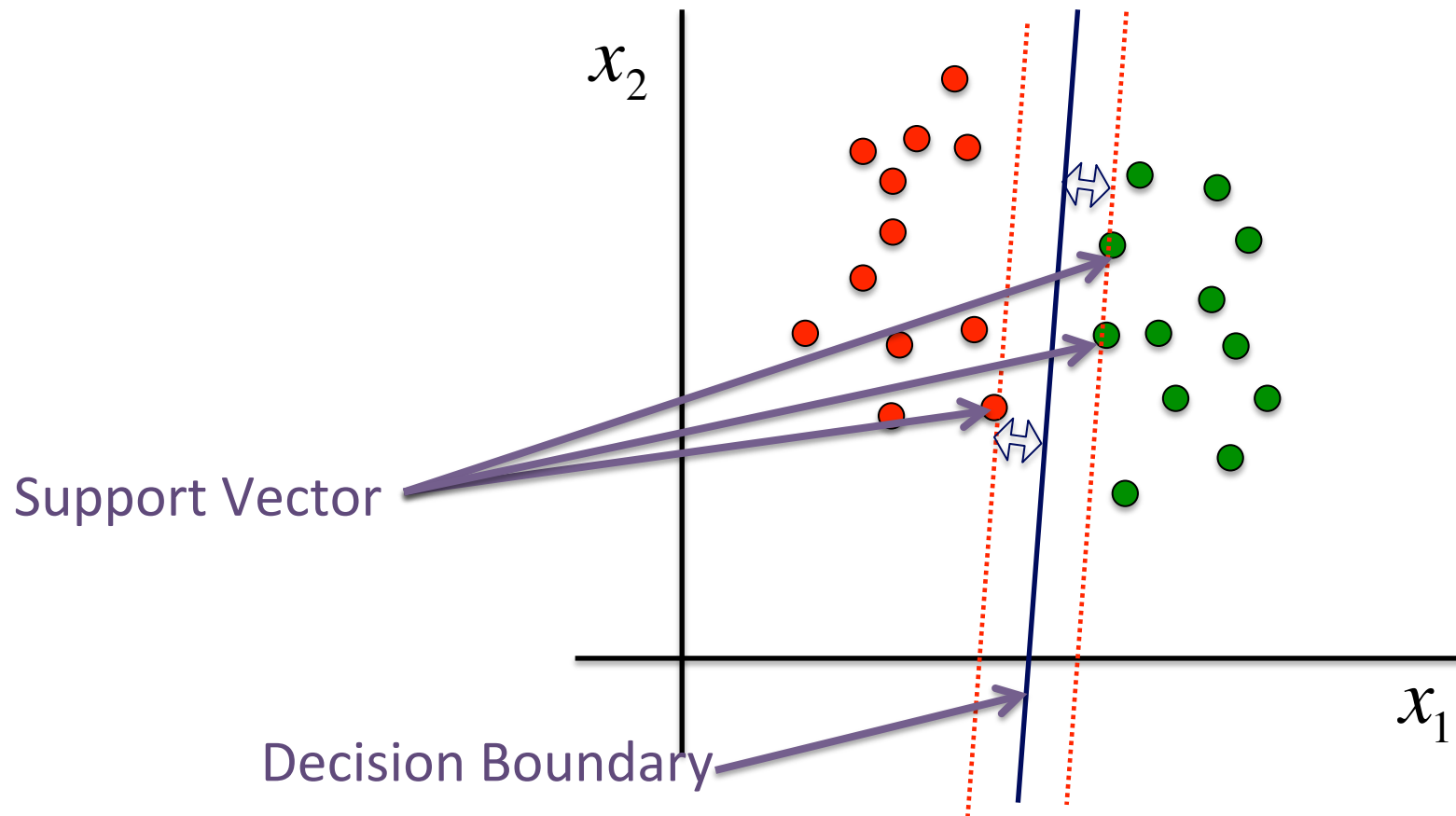
In the support vector machines the decision boundary is chosen to be the one for which the margin is maximized.

$$y(\bar{x}_n) = \bar{w}^T \phi(\bar{x})$$

$$\frac{t_n y(\bar{x}_n)}{\|\bar{w}\|}$$



In the support vector machines the decision boundary is chosen to be the one for which the margin is maximized.



An important property of support vector machines is that the determination of the model parameters corresponds to a convex optimization problem, and so any local solution is also a global optimum.

$$L(w, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n \left\{ t_n \bar{w} \phi(x_n) + b \right\} - 1$$

$$y(\bar{x}) = \sum_{n=1}^N a_n t_n k(\bar{x}, \bar{x}_n)$$

An important property of support vector machines is that the determination of the model parameters corresponds to a convex optimization problem, and so any local solution is also a global optimum.

Loss Function

$$L(w, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n \left\{ t_n \boxed{w \phi(x_n)} + b \right\} - 1$$

Lagrange Multiplier (points to a_n)
Ground truth (points to t_n)
Prediction (points to $w \phi(x_n)$)

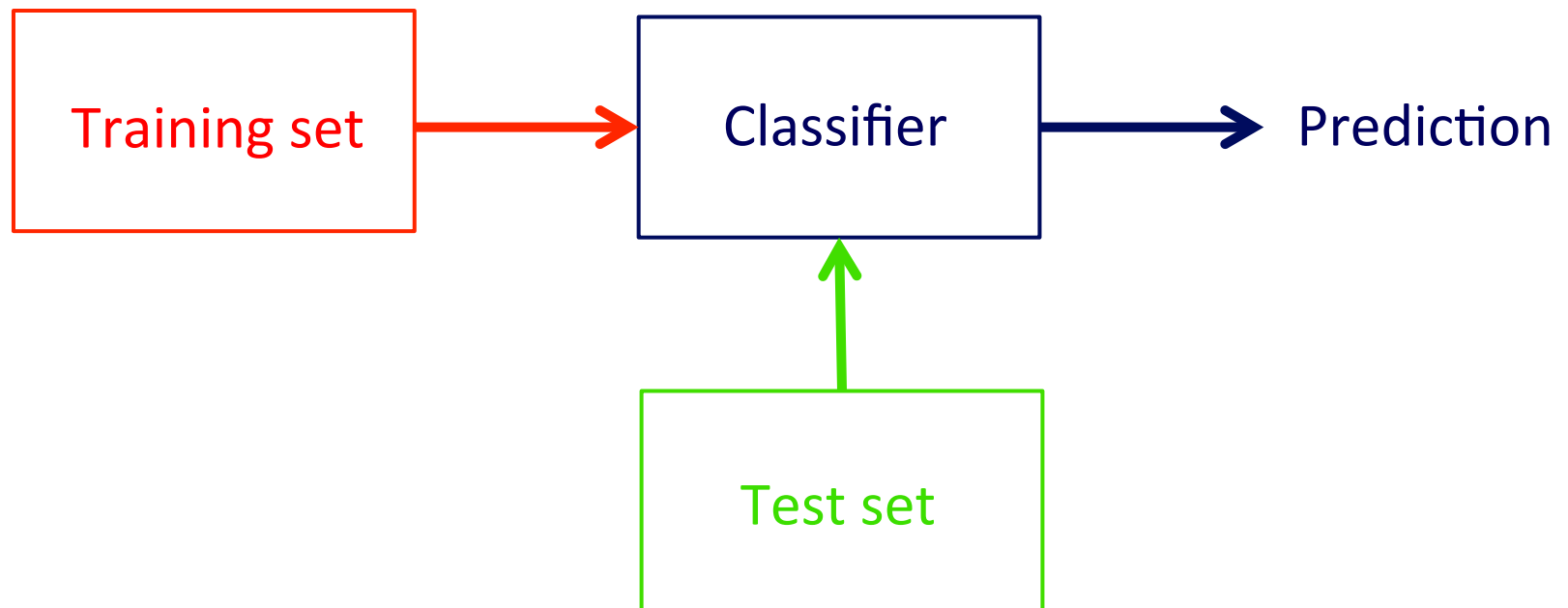
Model

$$y(\bar{x}) = \sum_{n=1}^N a_n t_n \boxed{k(\bar{x}, x_n)}$$

Kernel (points to $k(\bar{x}, x_n)$)

- The SVM is widely used in several domains.
- Several variants of the algorithm and kernel types have been proposed in the literature.
- The classifier, according also to the kernel used, requires the tuning of some hyper-parameters.

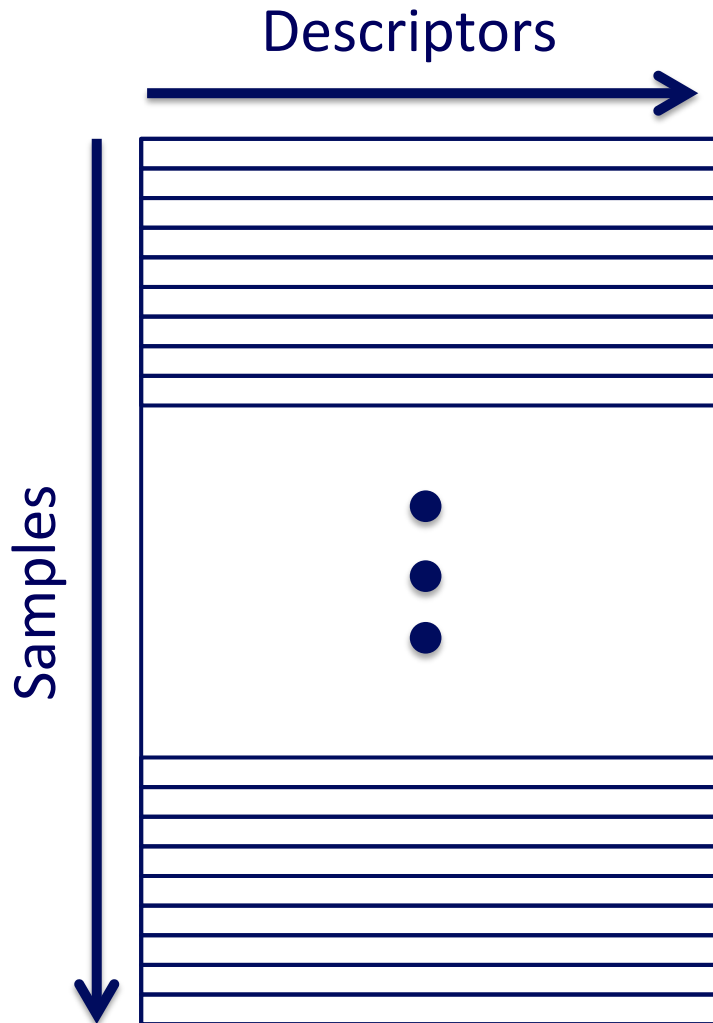
CLASSIFICATION SCHEME



The classification model depends on the samples used in the training stage. The ability of the system to recognize new samples depend on the parameters used.

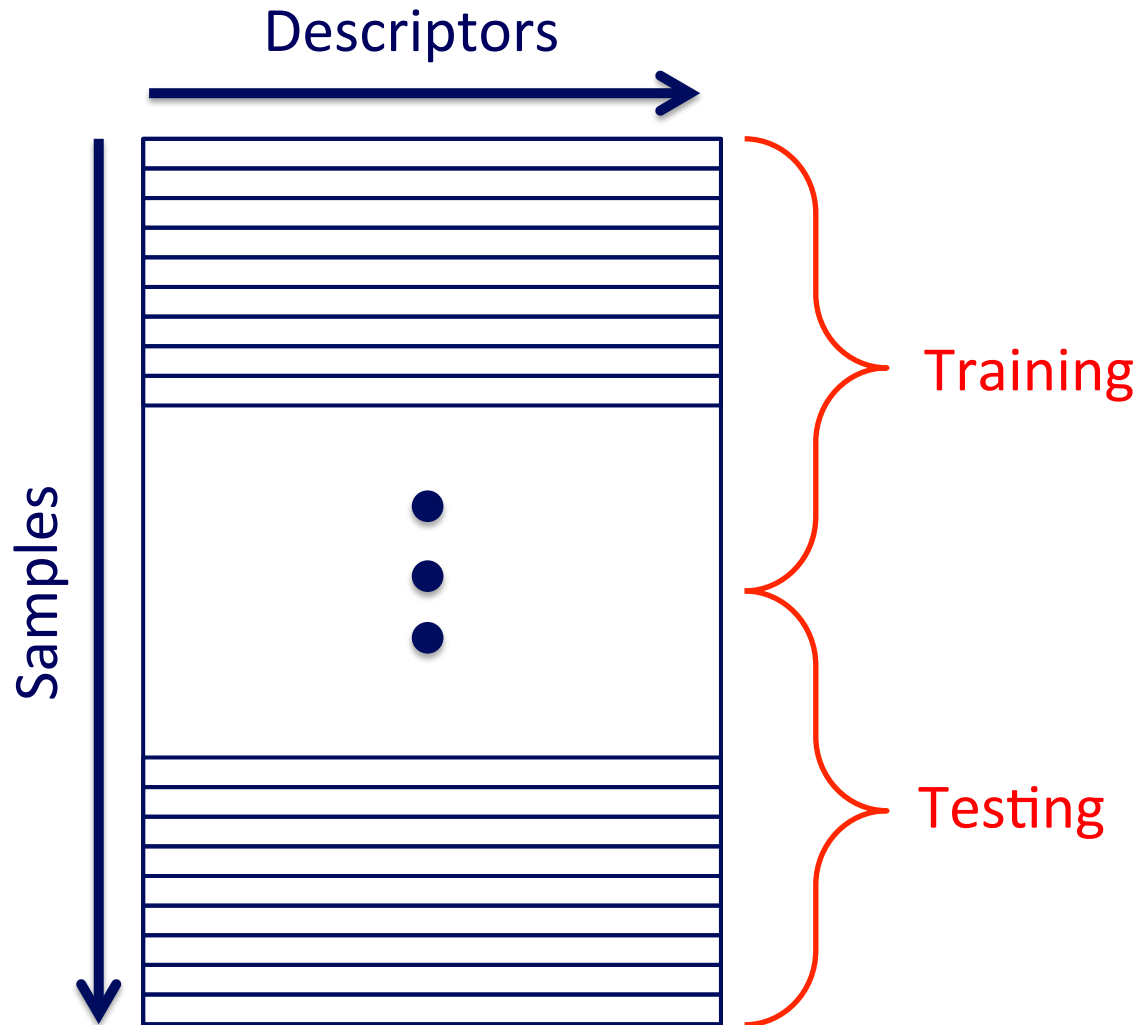
- How to perform tests independently of the samples used
- How to chose the optimal classifier' hyper-parameters

RANDOM SUB-SAMPLING

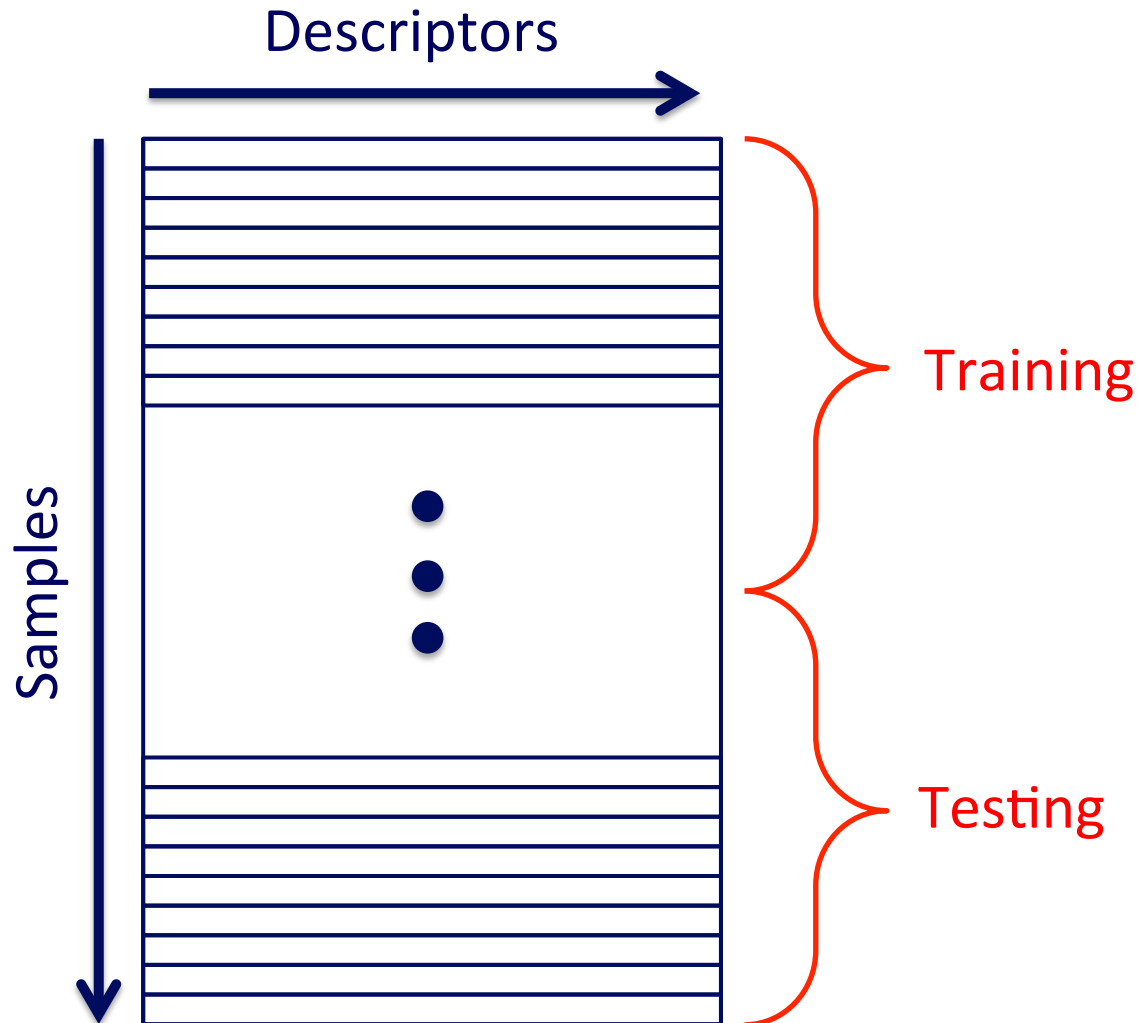


Which is the best way to generate the training set and the test set from the original dataset?

RANDOM SUB-SAMPLING

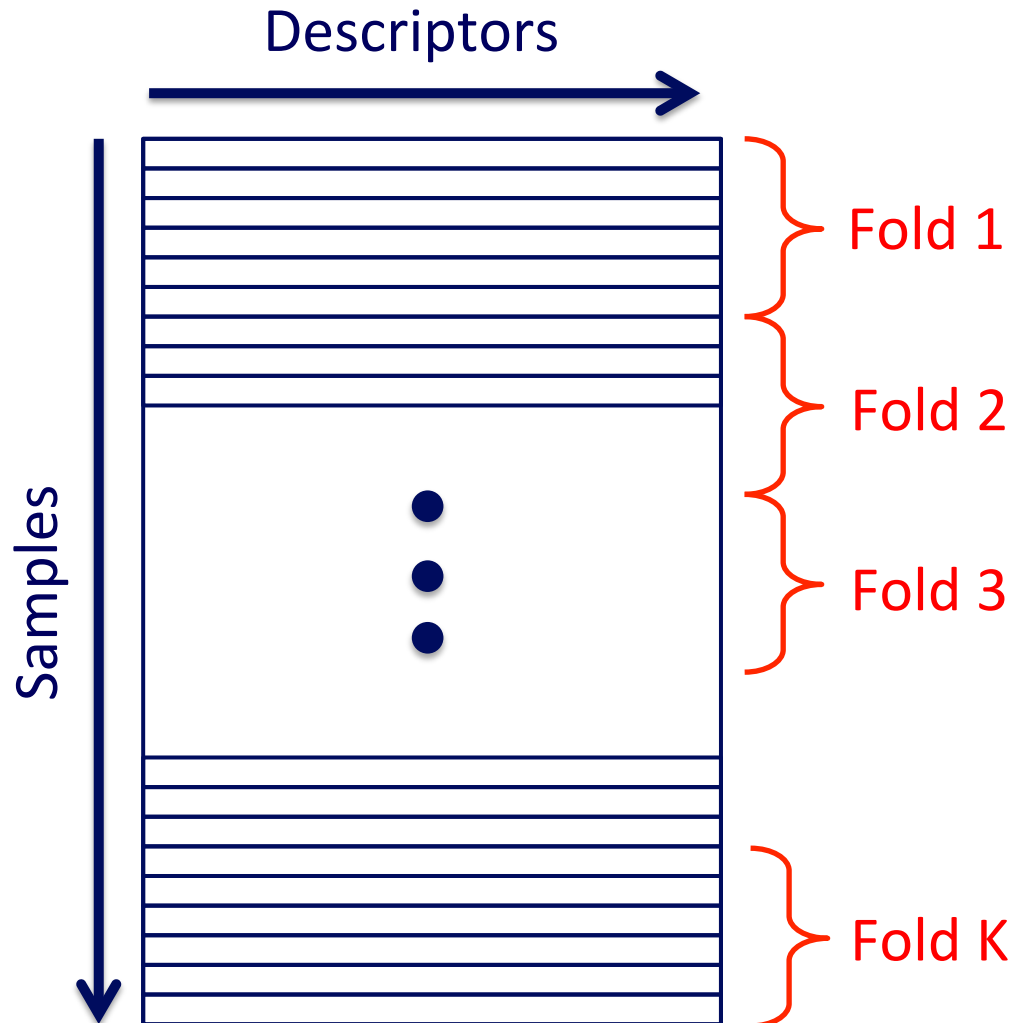


RANDOM SUB-SAMPLING



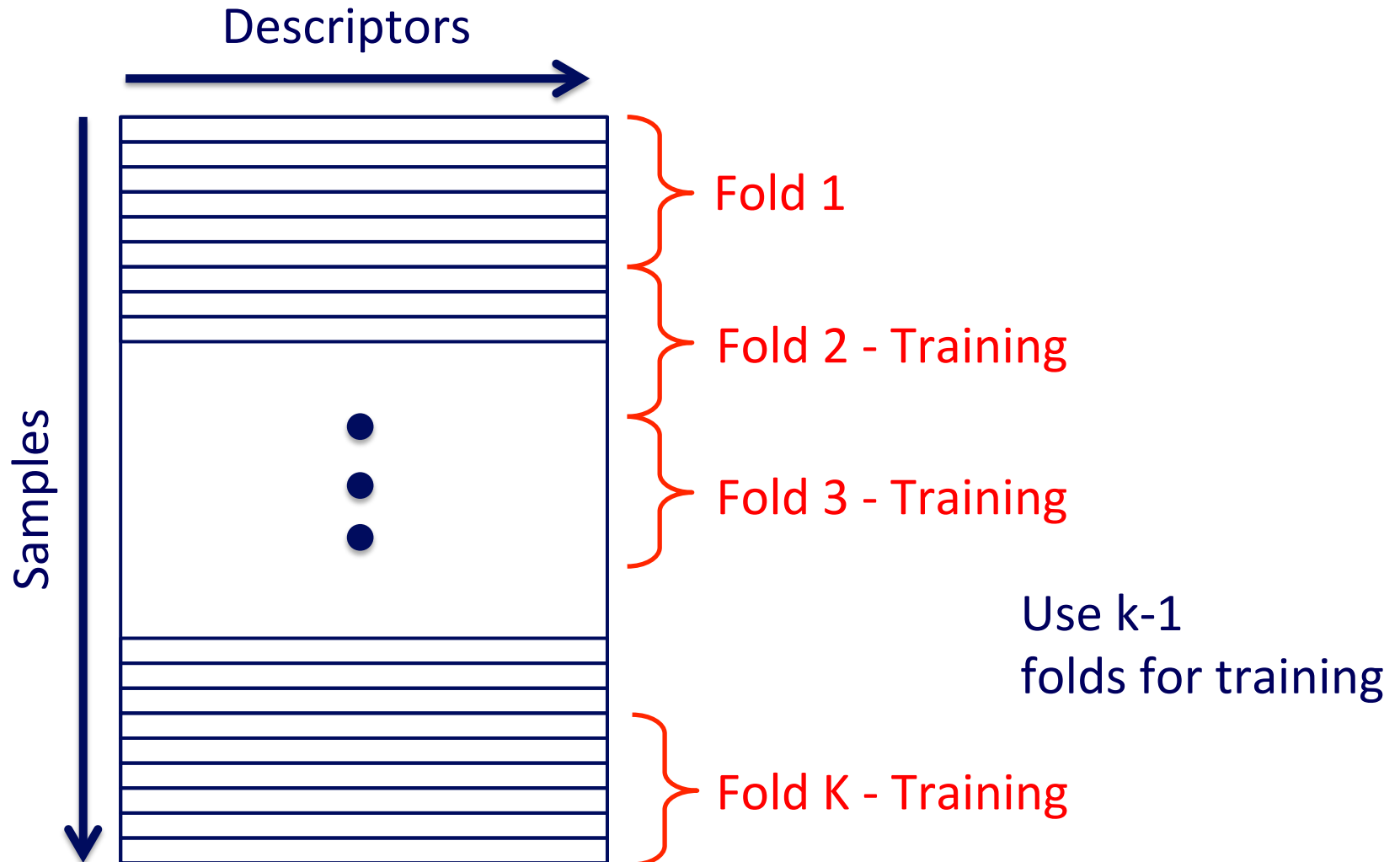
The model depends
on the samples used
in the training

K FOLD CROS-VALIDATION

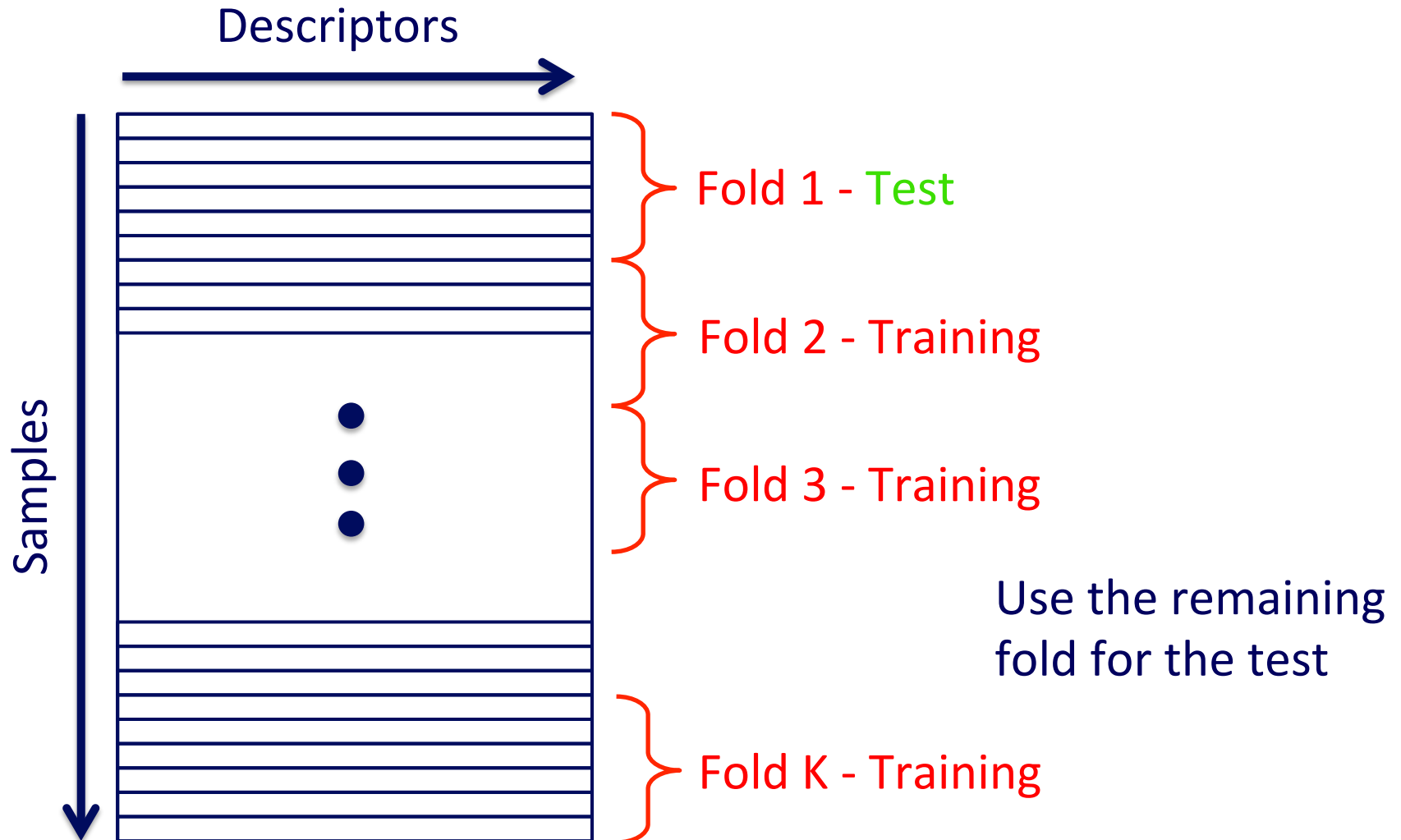


Define k
independent
folds

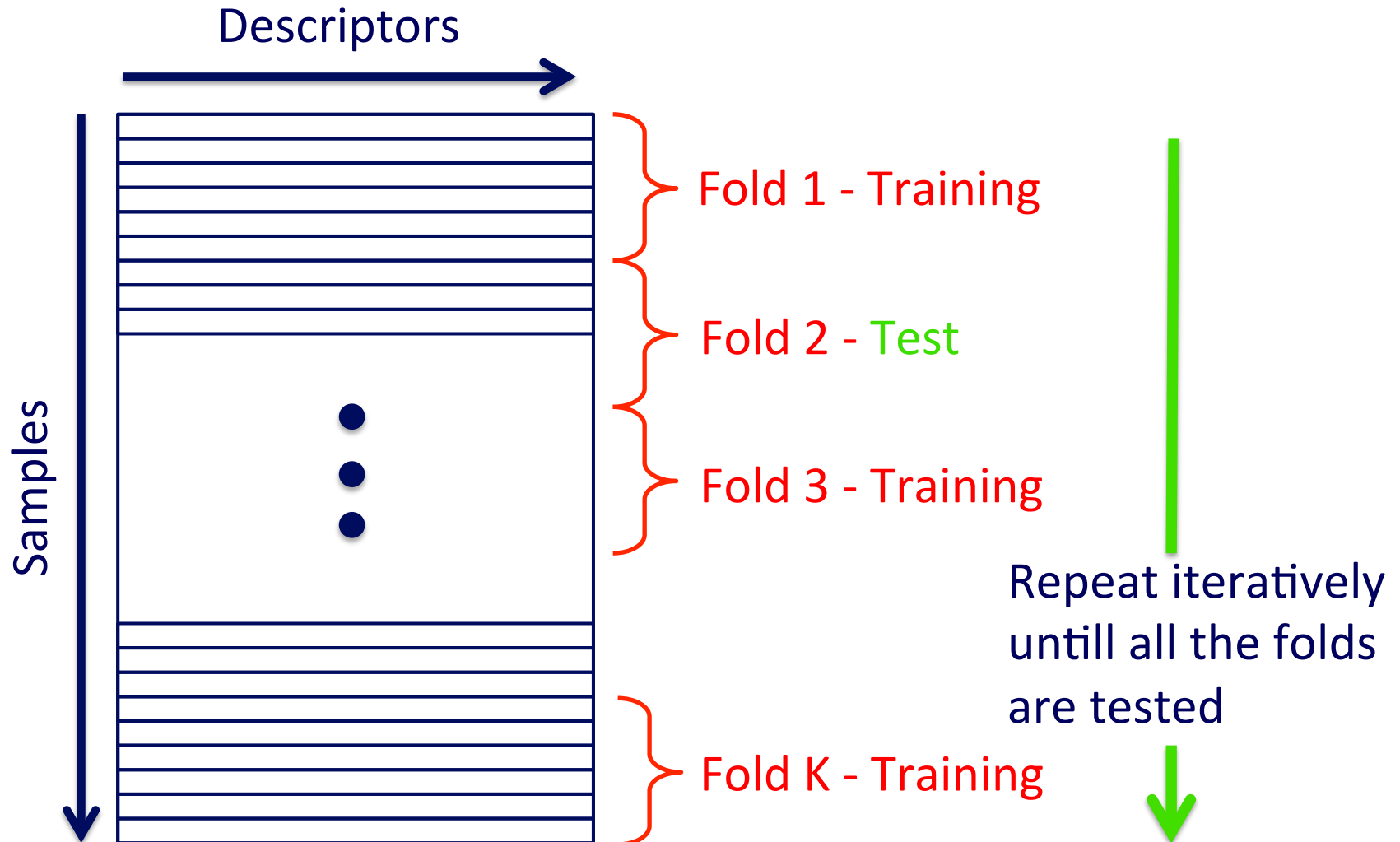
K FOLD CROS-VALIDATION



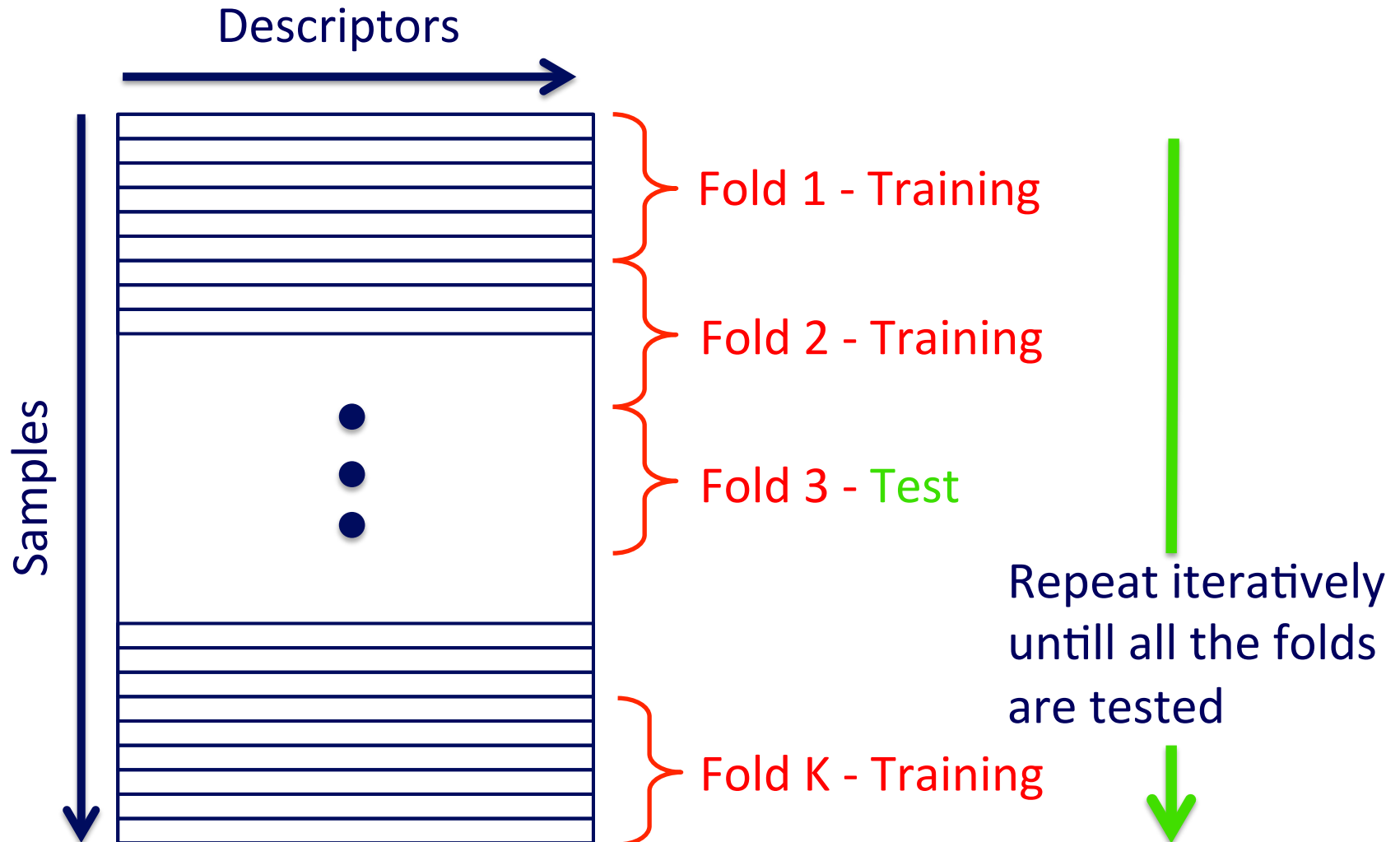
K FOLD CROS-VALIDATION



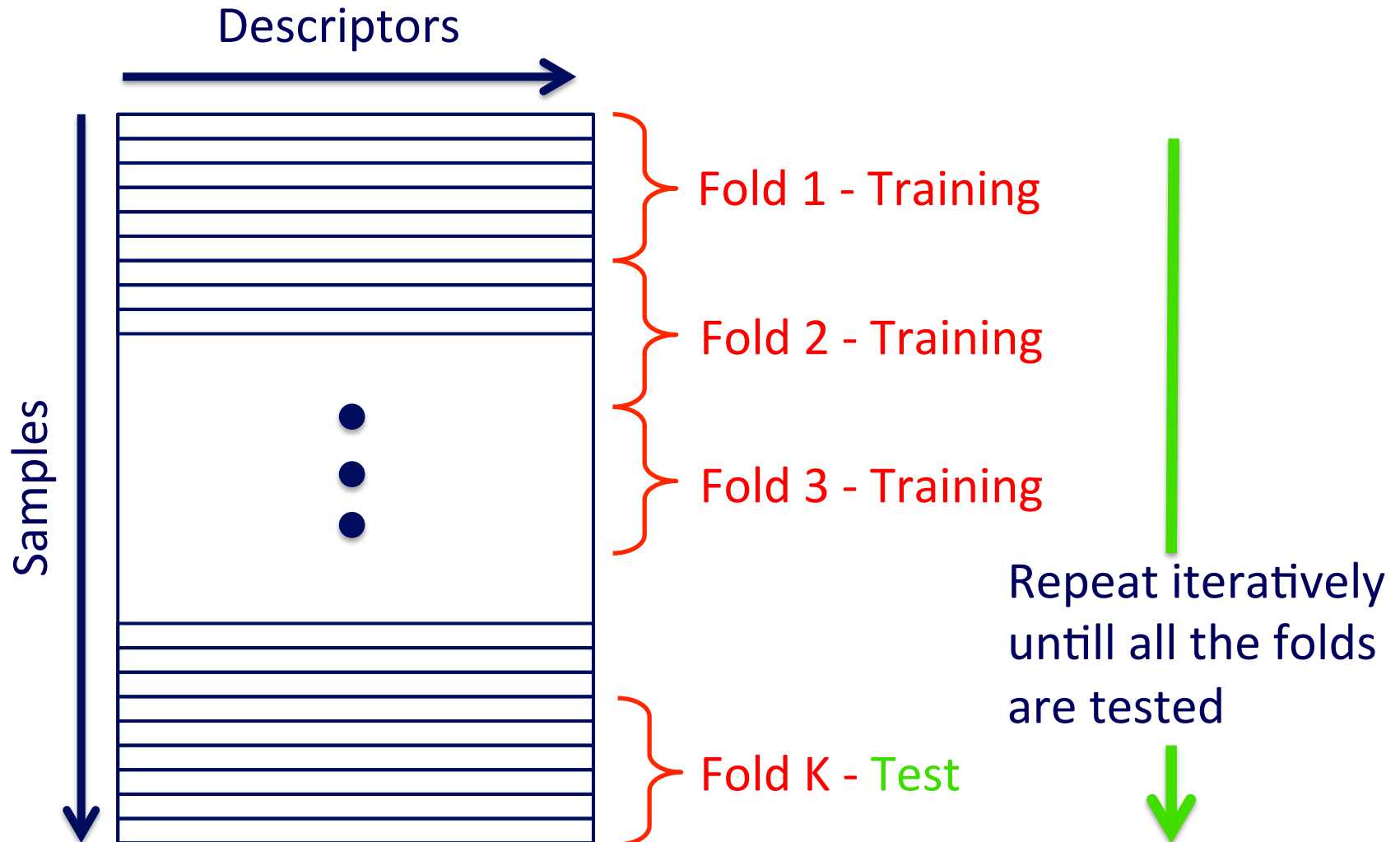
K FOLD CROS-VALIDATION



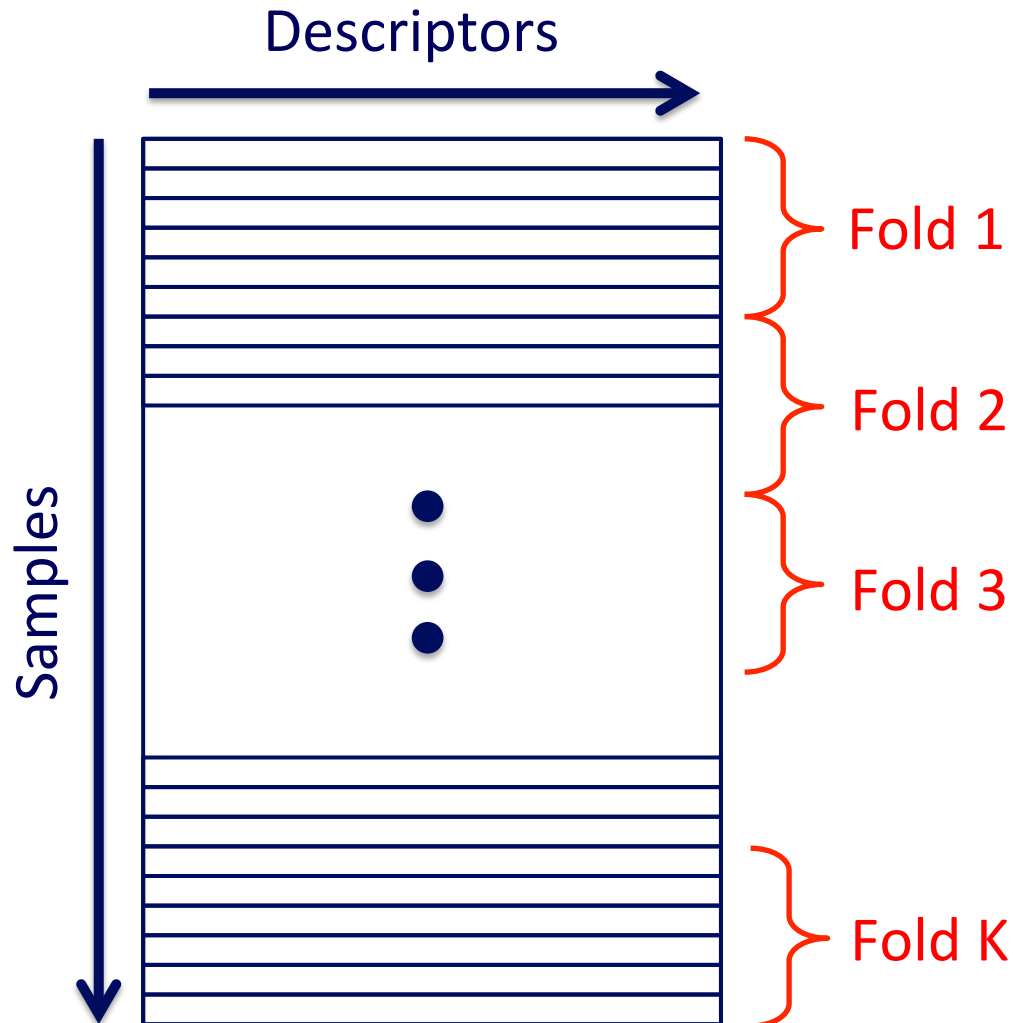
K FOLD CROS-VALIDATION



K FOLD CROS-VALIDATION



K FOLD CROS-VALIDATION




All the samples tested

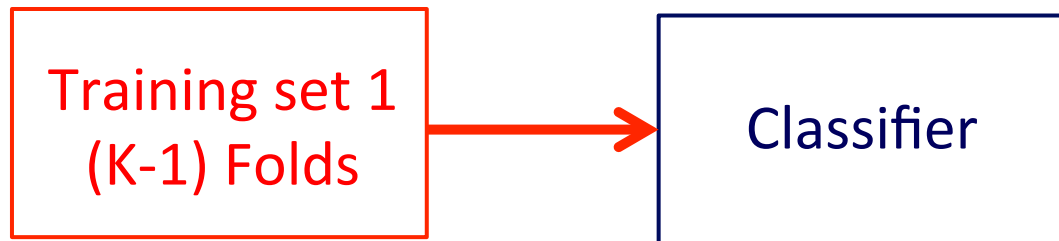


K FOLD CROS-VALIDATION

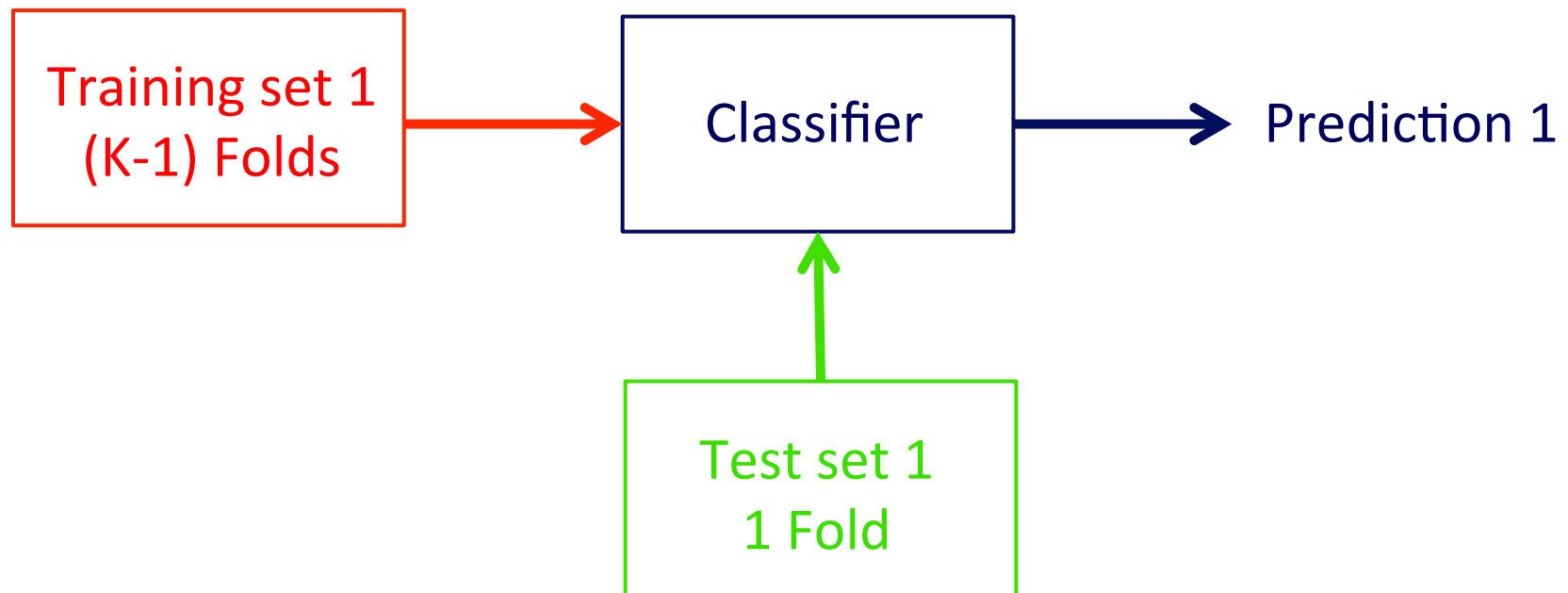
Training set 1
(K-1) Folds



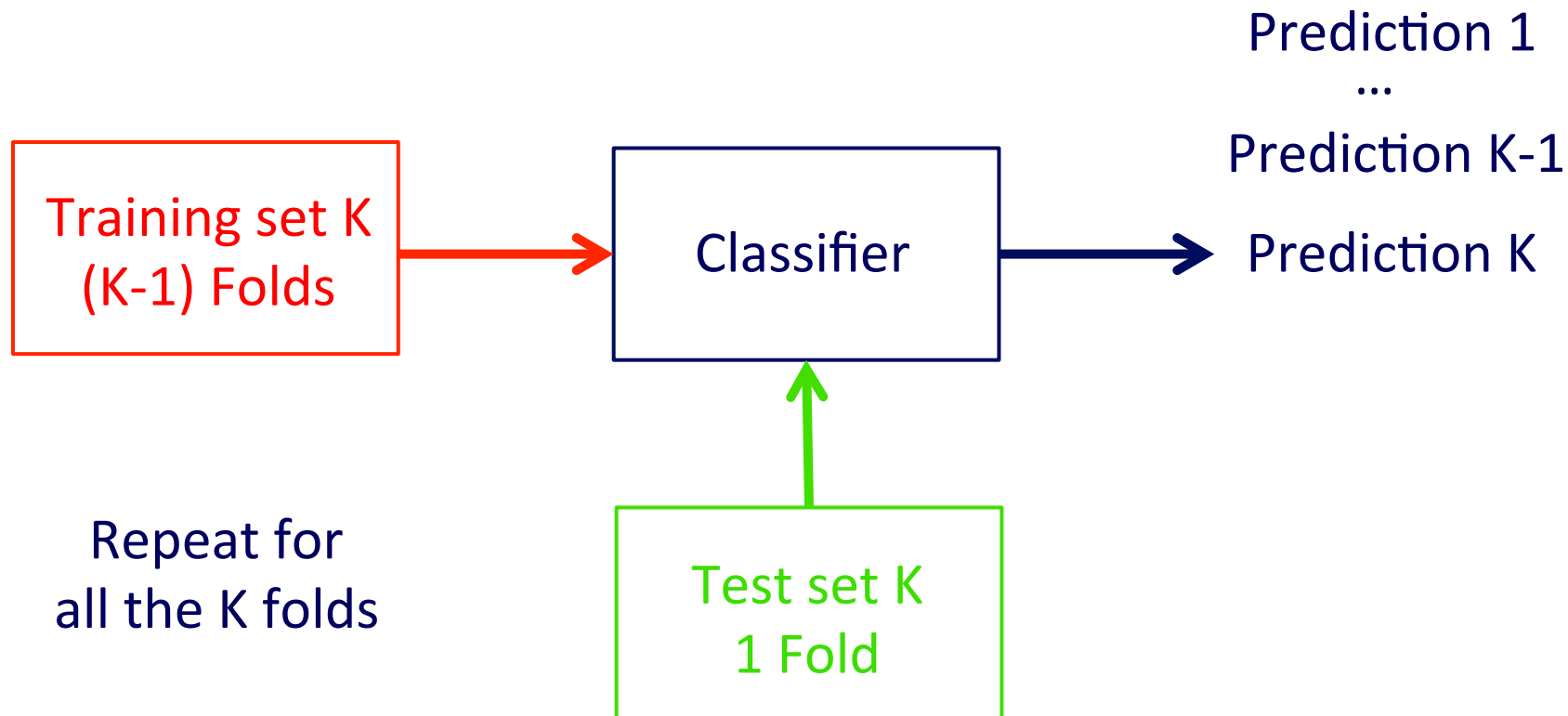
K FOLD CROS-VALIDATION

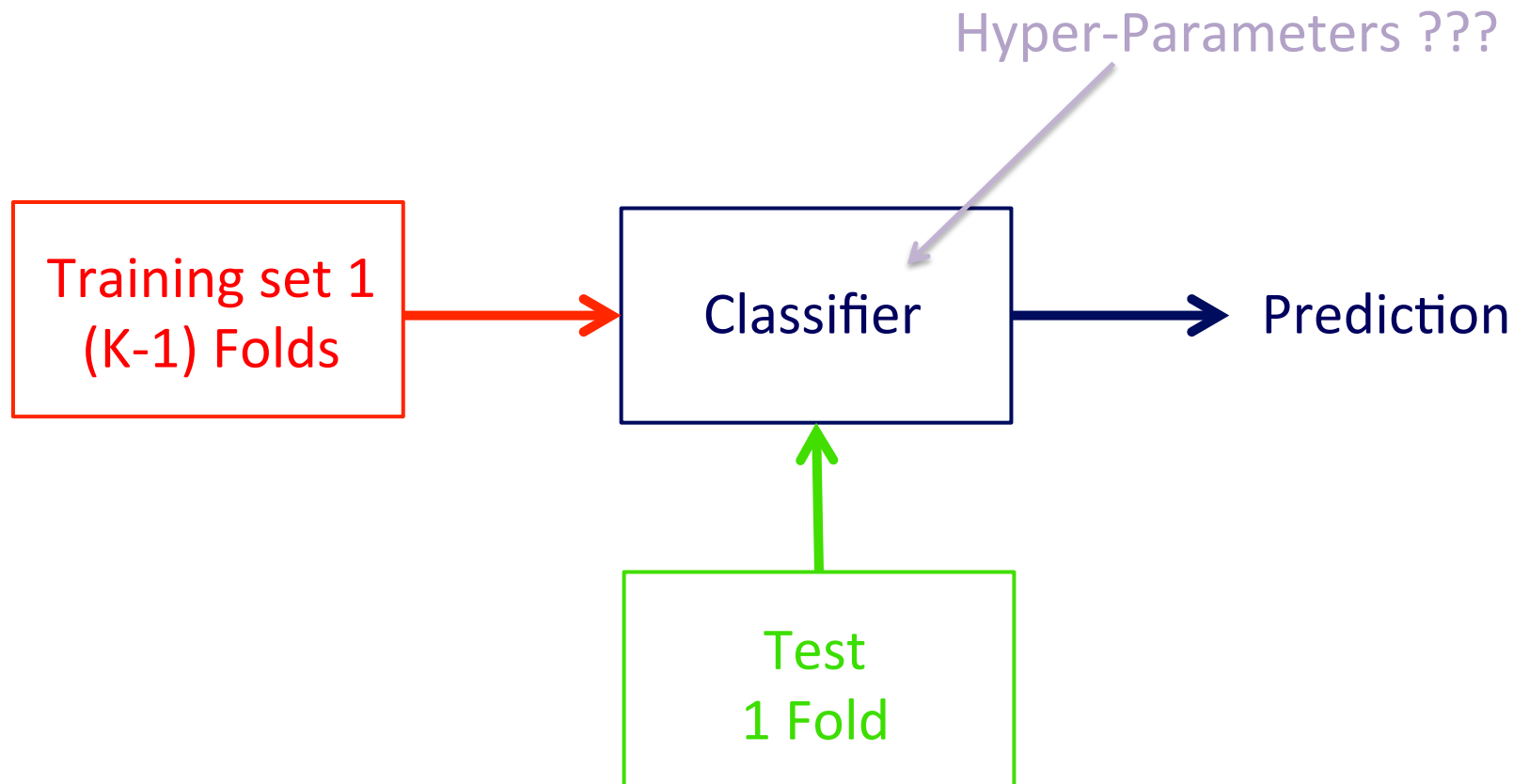


K FOLD CROS-VALIDATION



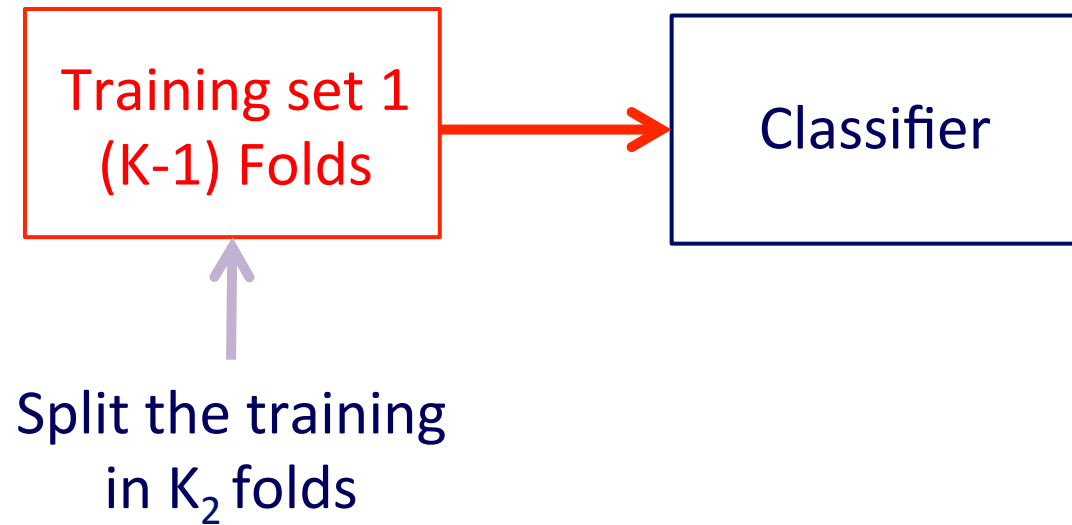
K FOLD CROSS-VALIDATION





Fix the value of the
Classifier parameters:

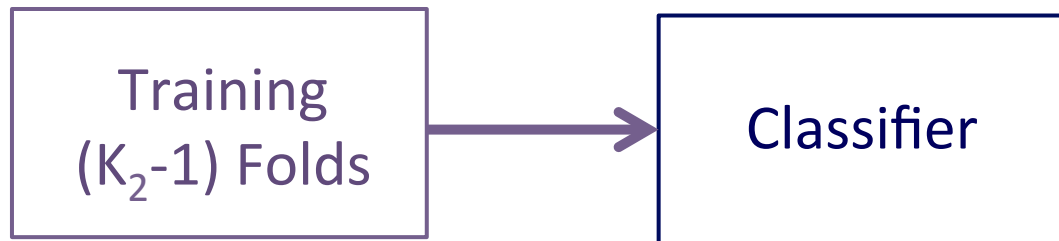
P_1, P_2, P_3, \dots



K FOLD CROS-VALIDATION
ON THE TRAINING SET

Fix the value of the
Classifier hyper-parameters:

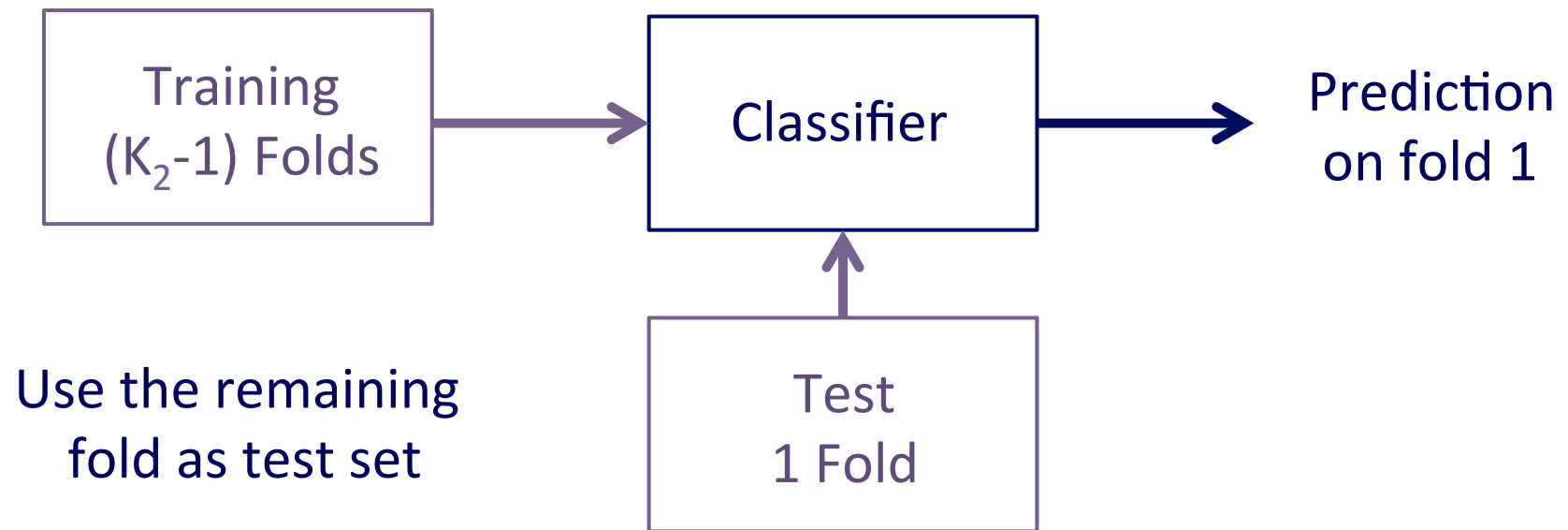
P_1, P_2, P_3, \dots



Use K_2 folds as
Training set

Fix the value of the
Classifier hyper-parameters:

P_1, P_2, P_3, \dots



**K FOLD CROS-VALIDATION
ON THE TRAINING SET**

Fix the value of the
Classifier hyper-parameters:
 P_1, P_2, P_3, \dots

Training
($K_2 - 1$) Fold

Classifier

Prediction
on fold 1

Prediction
on fold 2

⋮

Prediction
on fold K_2

Repeat for all
the K_2 folds

Test
1 Fold

**K FOLD CROS-VALIDATION
ON THE TRAINING SET**

Fix the value of the
Classifier hyper-parameters:
 P_1, P_2, P_3, \dots

Training
($K_2 - 1$) Fold

Classifier

Prediction
on fold 1

Prediction
on fold 2

⋮

Prediction
on fold K_2

Repeat for all
the K_2 folds

Test
1 Fold

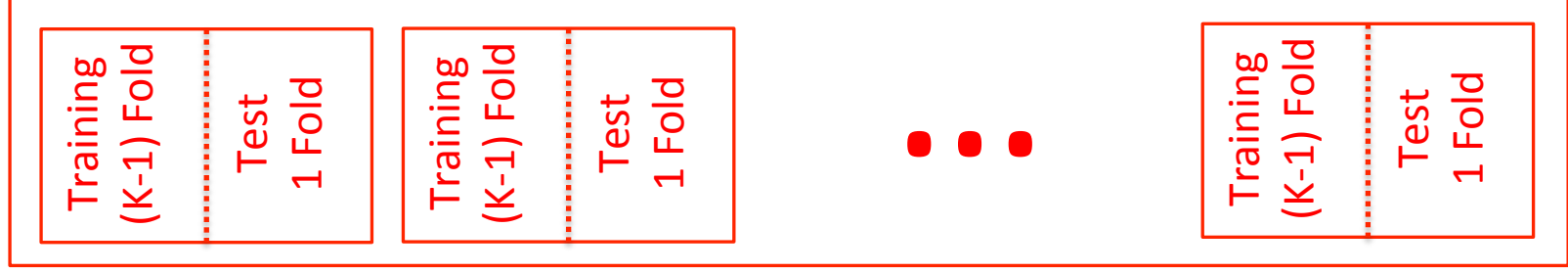
Average performance of the system using hyper-parameters:
 P_1, P_2, P_3, \dots



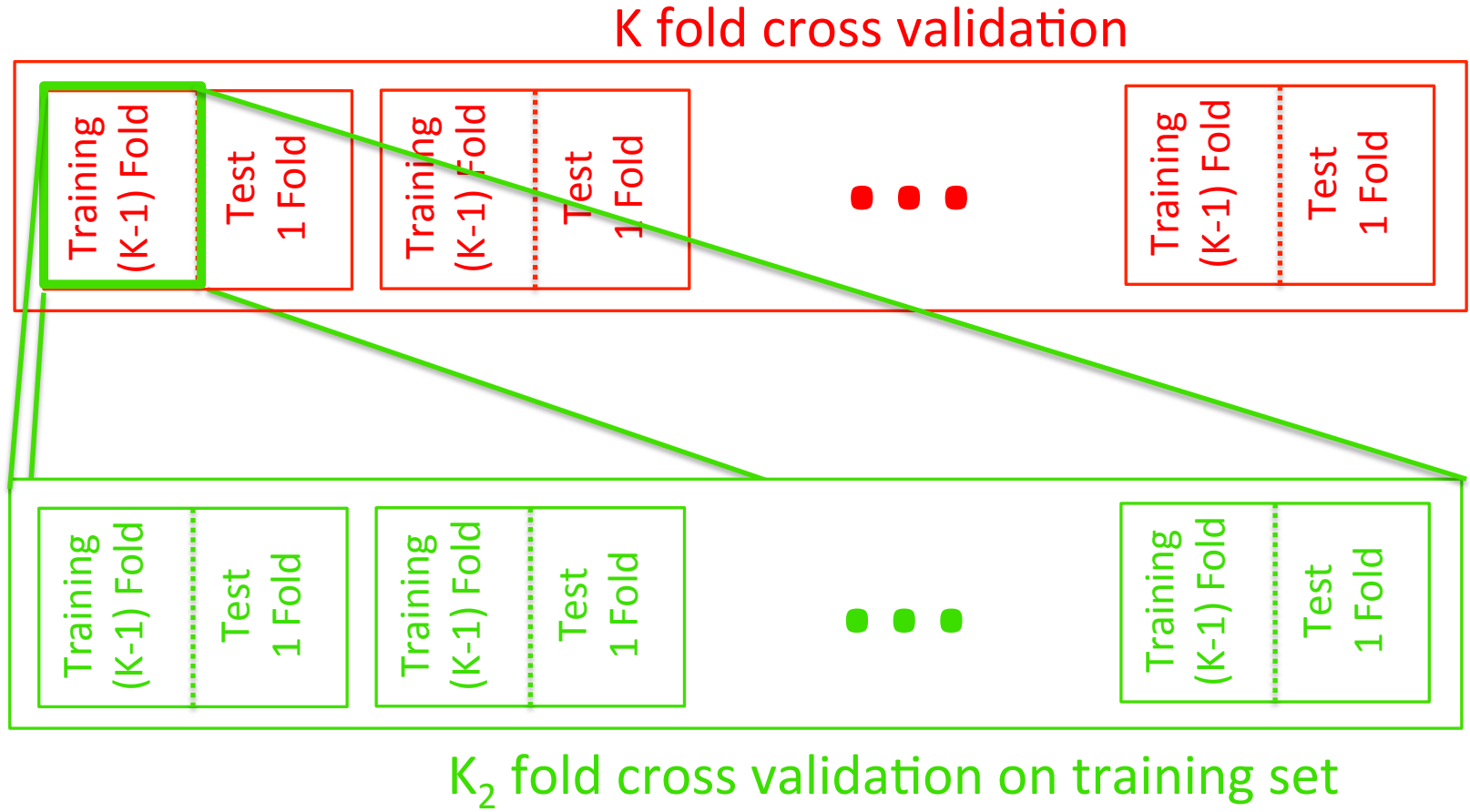
K FOLD CROS-VALIDATION
ON THE TRAINING SET

- Repeat the test with several hyper-parameters values
- Chose the set of hyper-parameters that provides the best performance
- Train the classifier using the chosen hyper-parameters

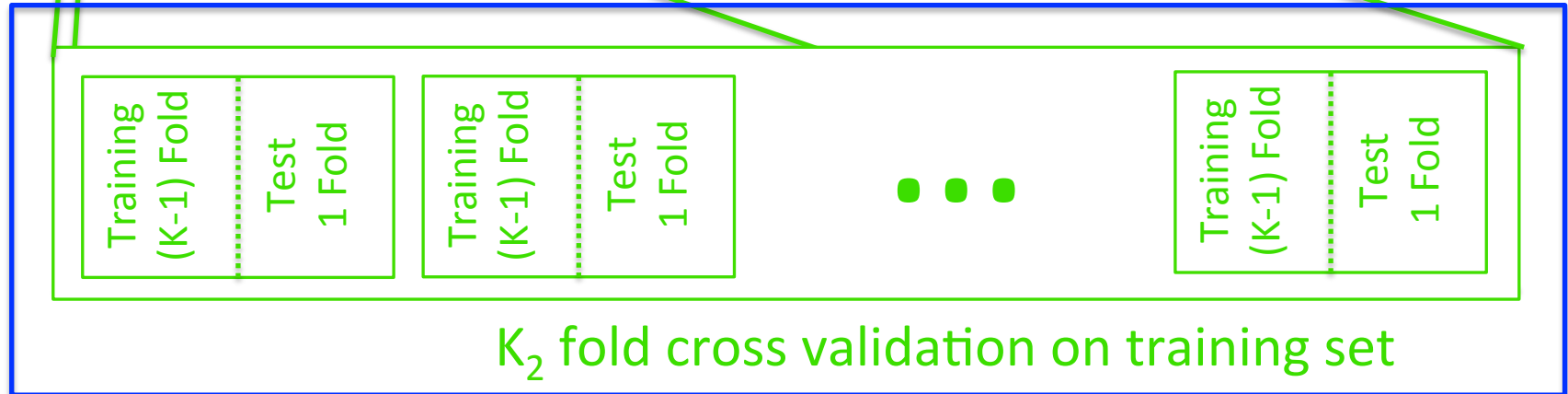
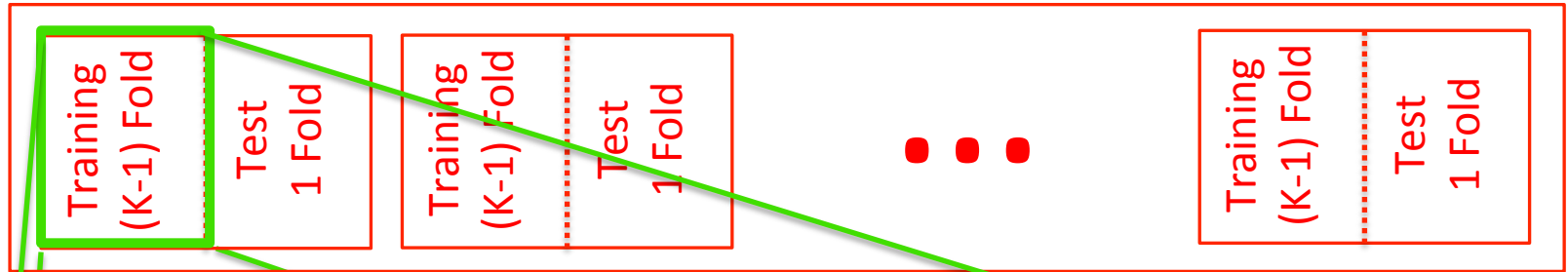
K fold cross validation



HOW TO PERFORM EXPERIMENTS

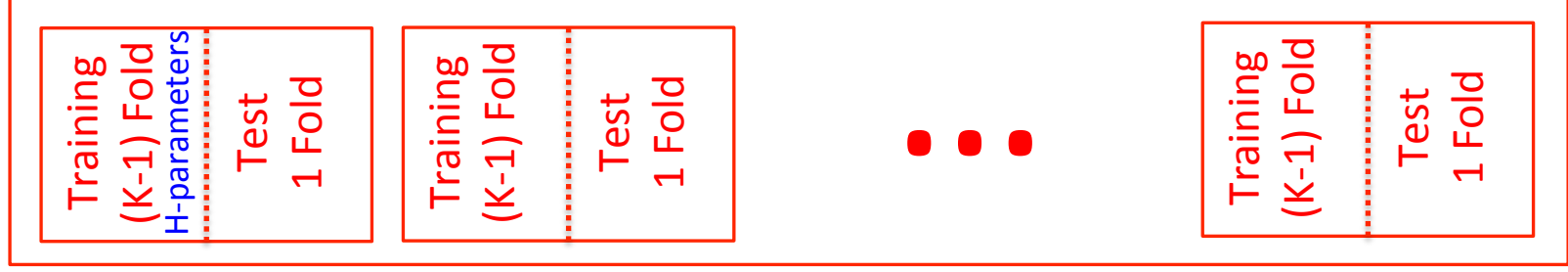


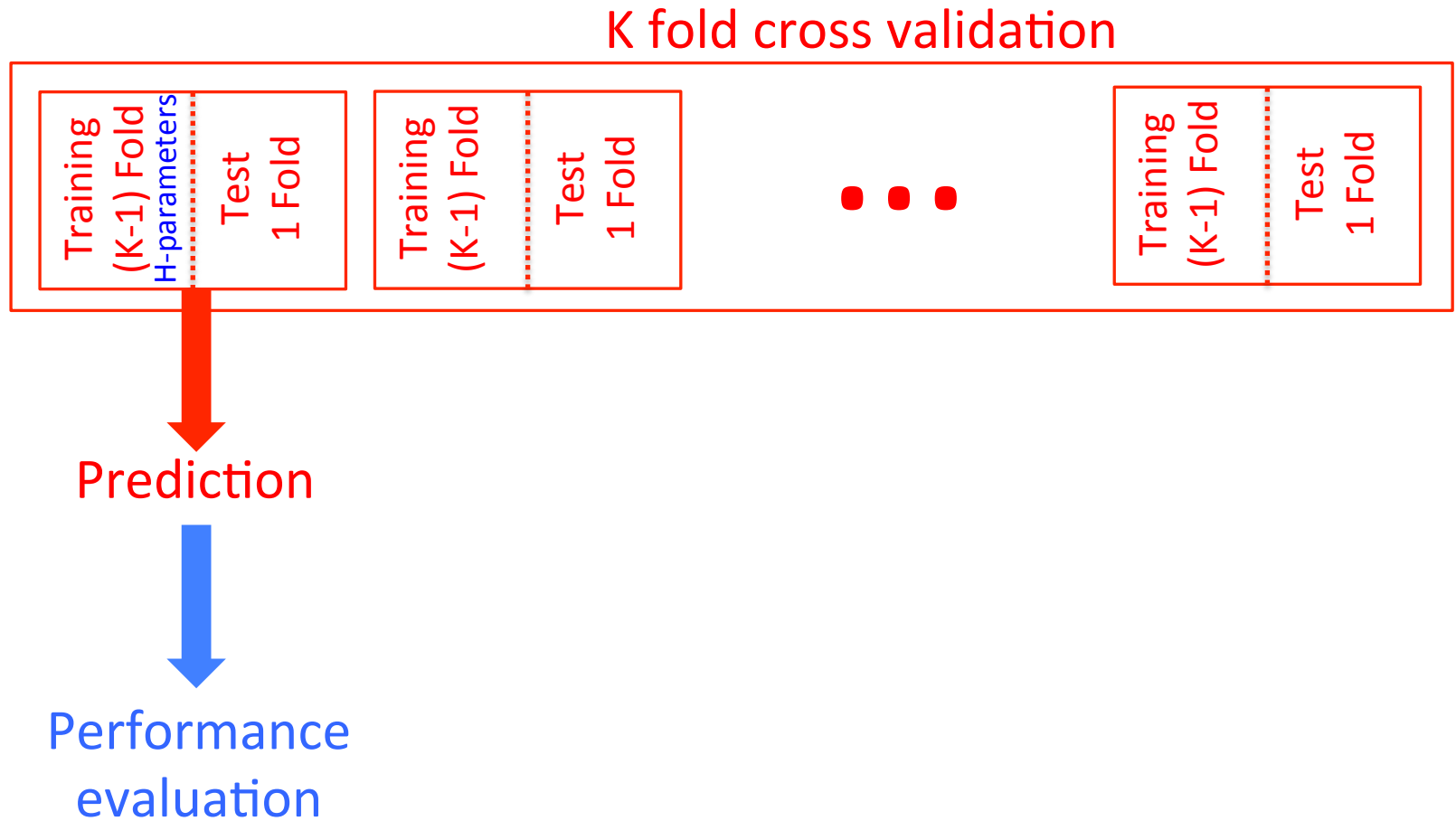
K fold cross validation

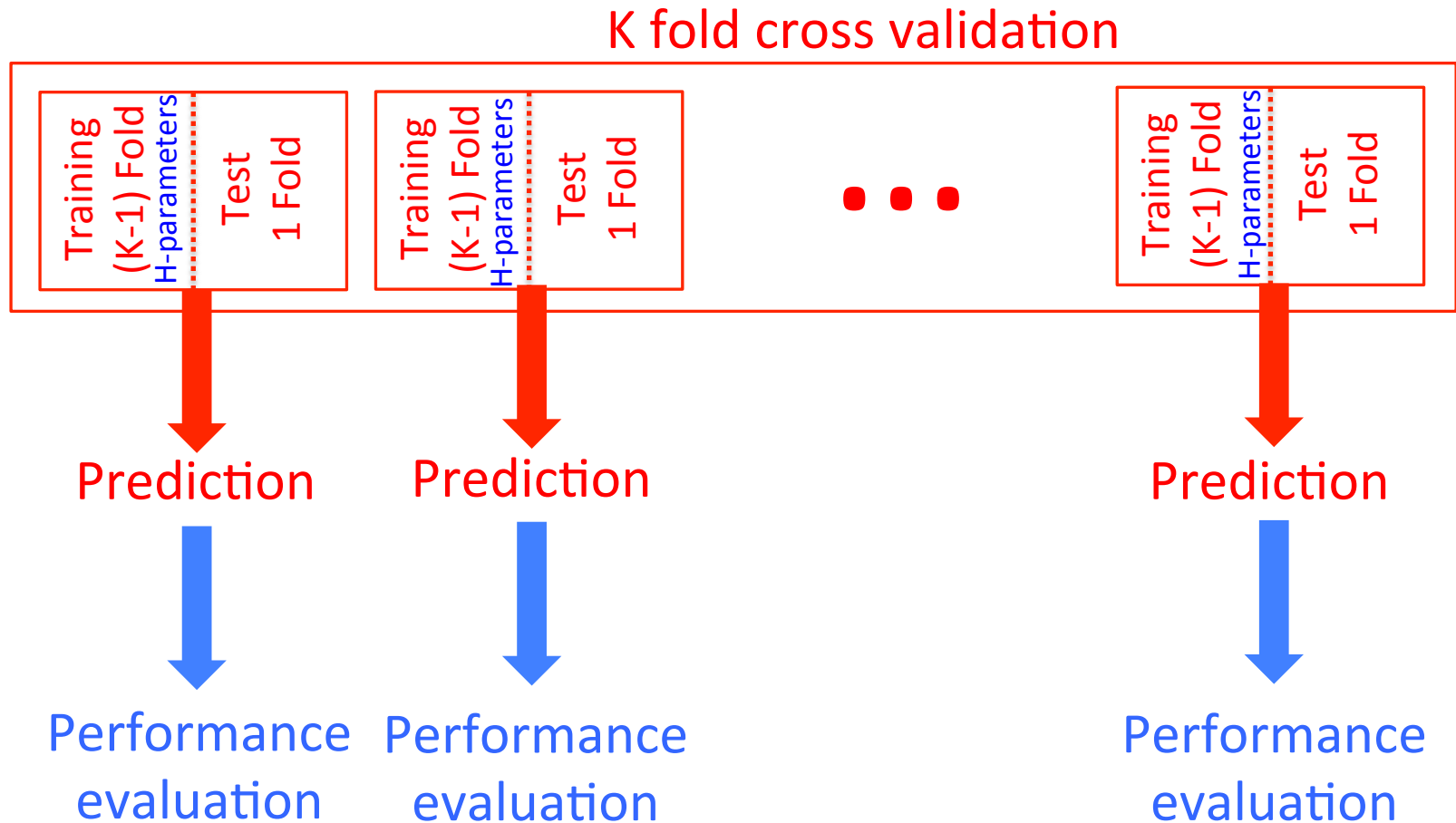


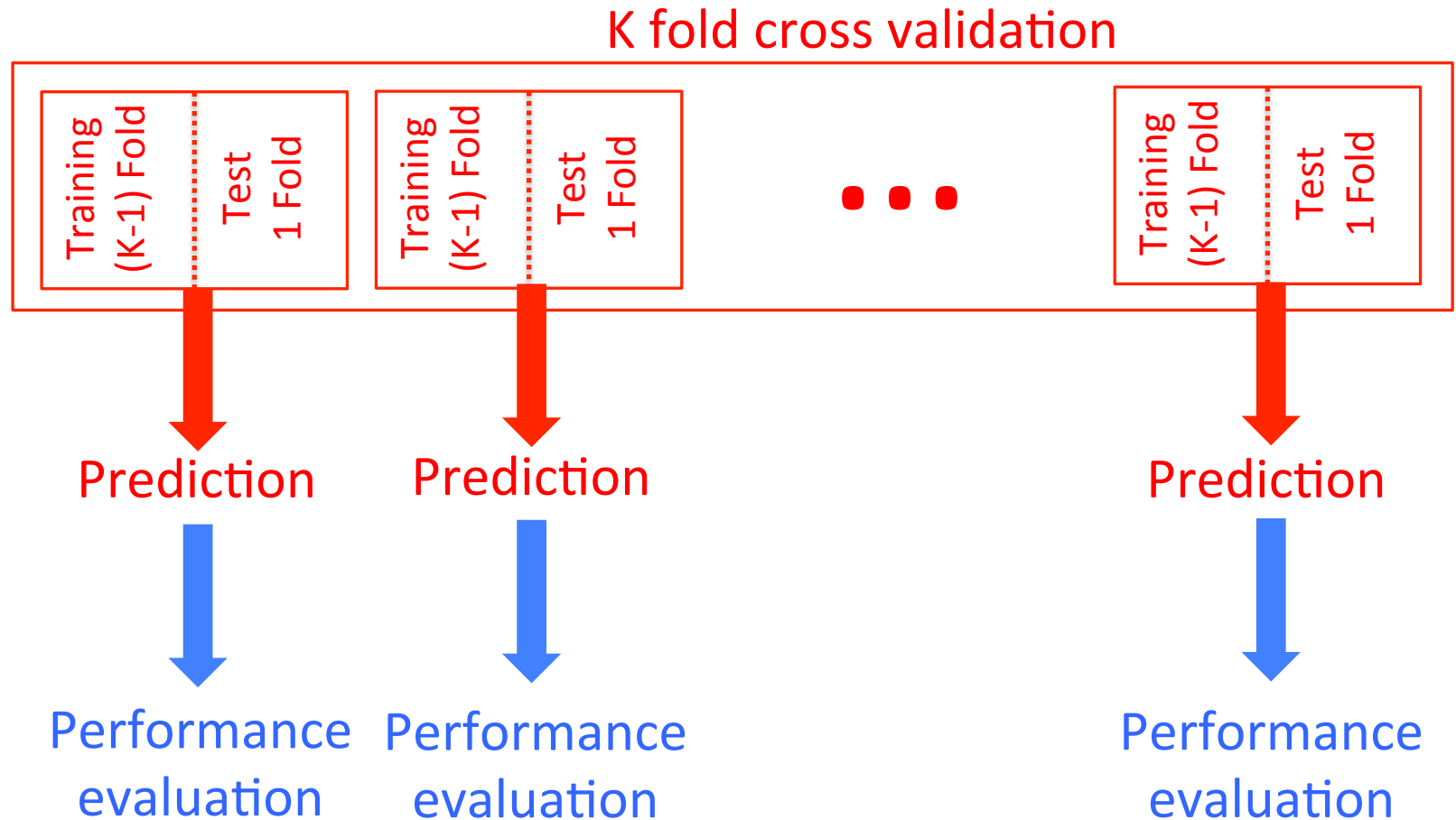
Repeat for each set of hyper-parameters

K fold cross validation





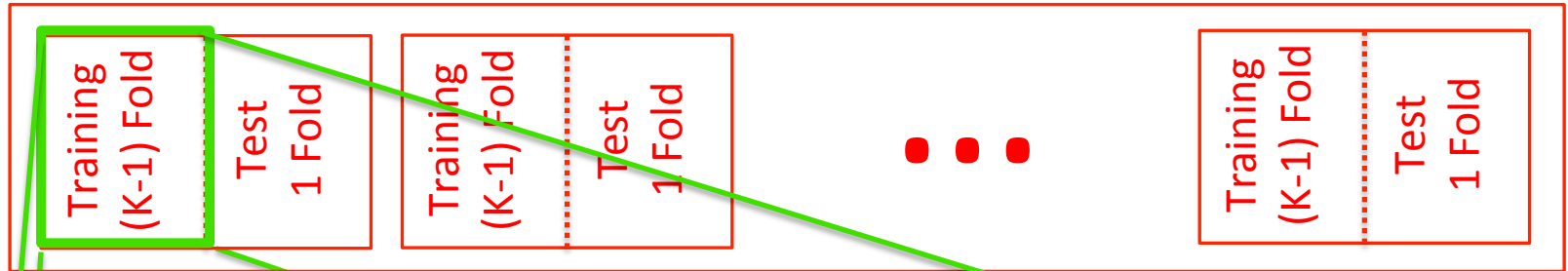




System performance are computed as average of performance among the folds

STRATIFIED CV

K fold cross validation



K₂ fold cross validation on training set

Repeat for each set of parameters

Repeat for each layer of the stratification

CONFUSION MATRIX

		Ground Truth	
		p	n
Predicted Class	\hat{p}	True Positive	False Positive
	\hat{n}	False Negative	True Negative
		P	N

ACCURACY

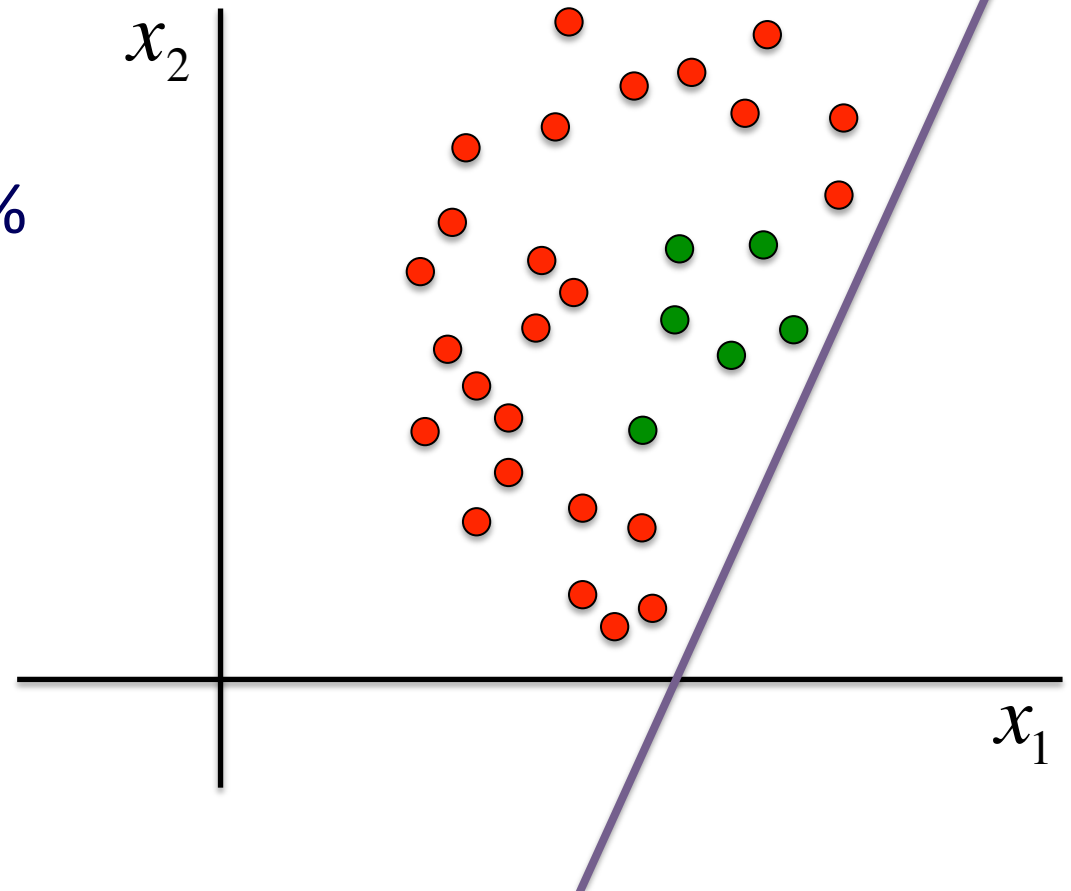
$$\text{Acc} = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}}$$

The number of samples correctly classified over the total number of samples

	p	n
\hat{p}	True Positive	False Positive
\hat{n}	False Negative	True Negative
	P	N

ACCURACY LIMITATIONS

Accuracy = 80%



$$\text{Recall} = \frac{TP}{P}$$

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

\hat{p}

\hat{N}

	p	n
\hat{p}	True Positive	False Positive
\hat{N}	False Negative	True Negative
	P	N

OTHER METRICS

$$\text{G-mean} = \left(\frac{TP}{P} \cdot \frac{TN}{N} \right)^{\frac{1}{2}} \quad \tilde{P}$$

$$\text{Class acc} = \frac{\frac{TP}{P} + \frac{TN}{N}}{2} \quad \tilde{N}$$

$$\text{F-measure} = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

	p	n
P	True Positive	False Positive
N	False Negative	True Negative

- Pattern Recognition and Machine Learning; C.M.Bishop
- Pattern Classification; R.O. Duda, P.E. Hart, D.G. Stork
- Machine Learning; T.M.Mitchell
- Statistics and Data with R; Y. Cohen, J.Y. Cohen

