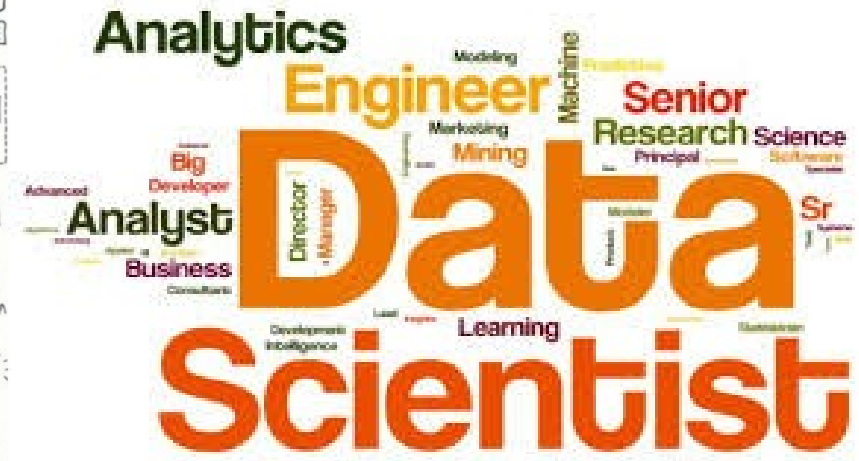
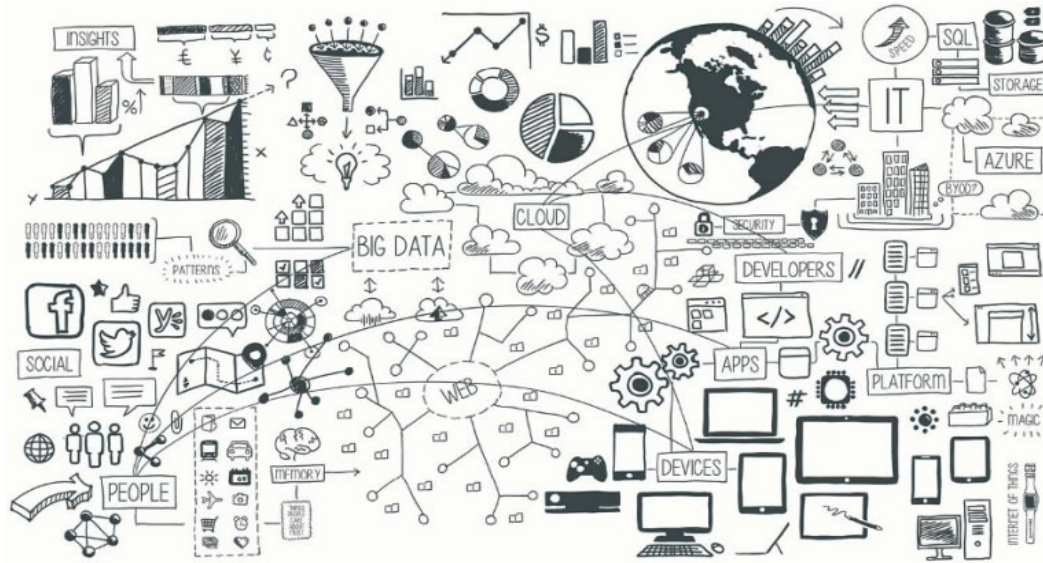


Data Mining (Minería de Datos)

Evaluación, sobreajuste y validación cruzada (cross-validation)



Sixto Herrera

Grupo de Meteorología
Univ. de Cantabria – CSIC
MACC / IFCA

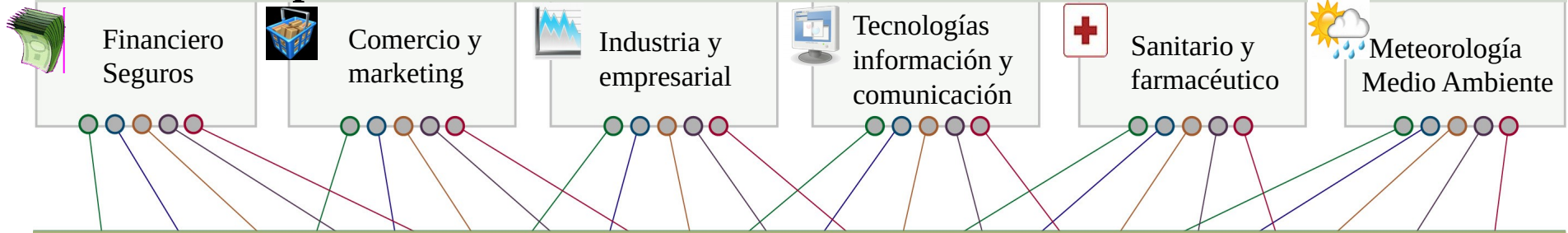


NOTA: Las líneas de código de R en esta presentación se muestran sobre un fondo gris.

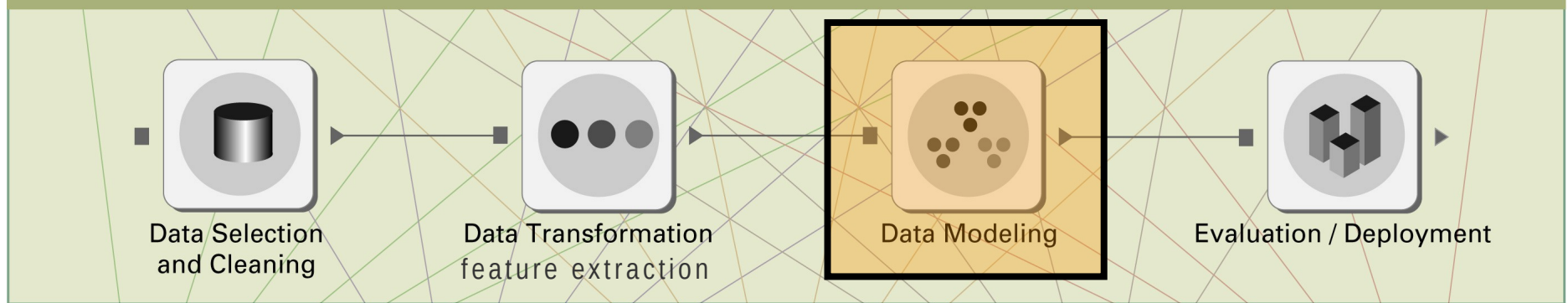
Cross-Validation

Oct	29	Aplazada (sesión de refuerzo)
	31	Presentación, introducción y perspectiva histórica
Nov	5	Paradigmas, problemas canonicos y data challenges
	7	Reglas de asociación
	12	Práctica: Reglas de asociación
	14	Evaluación, sobrejuste y crossvalidacion
	19	Práctica: Cross-validación
	21	Arboles de clasificacion y decision
	26	Practica: Árboles de clasificación
		T01. Datos discretos
	28	Técnicas de vecinos cercano (k-NN)
Dic	3	Práctica: Vecinos cercanos
	5	Reducción de dimensión no lineal
	10	Práctica: Reducción de dimensión no lineal
		T02. Clasificación
	12	Árboles de clasificación y regresion (CART)
	17	Práctica: Árboles de clasificación y regresion (CART)
	19	Ensembles: Bagging and Boosting
Ene	7	Práctica Random Forests
	9	Práctica Gradient boosting
		T03. Prediccion
	14	Técnicas de agrupamiento
	16	Practica: Técnicas de agrupamiento
	21a	Practica: El paquete CARET
	21b	Examen

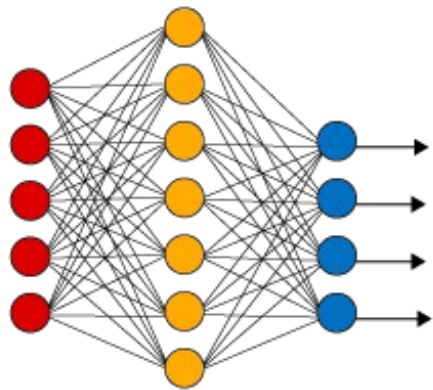
Sectores de aplicación



Proceso de Minería de Datos

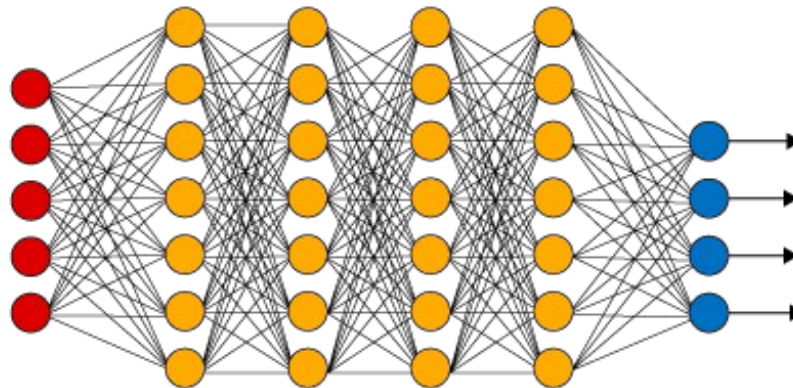


Simple Neural Network



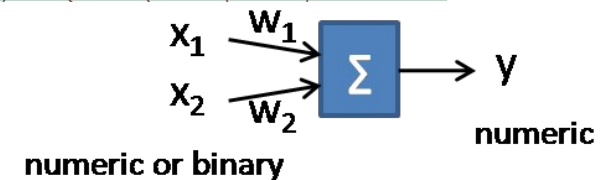
Input Layer

Deep Learning Neural Network



Hidden Layer

Output Layer



$$y = w_0 + w_1x_1 + w_2x_2$$

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{x}^T \cdot \mathbf{w}$$

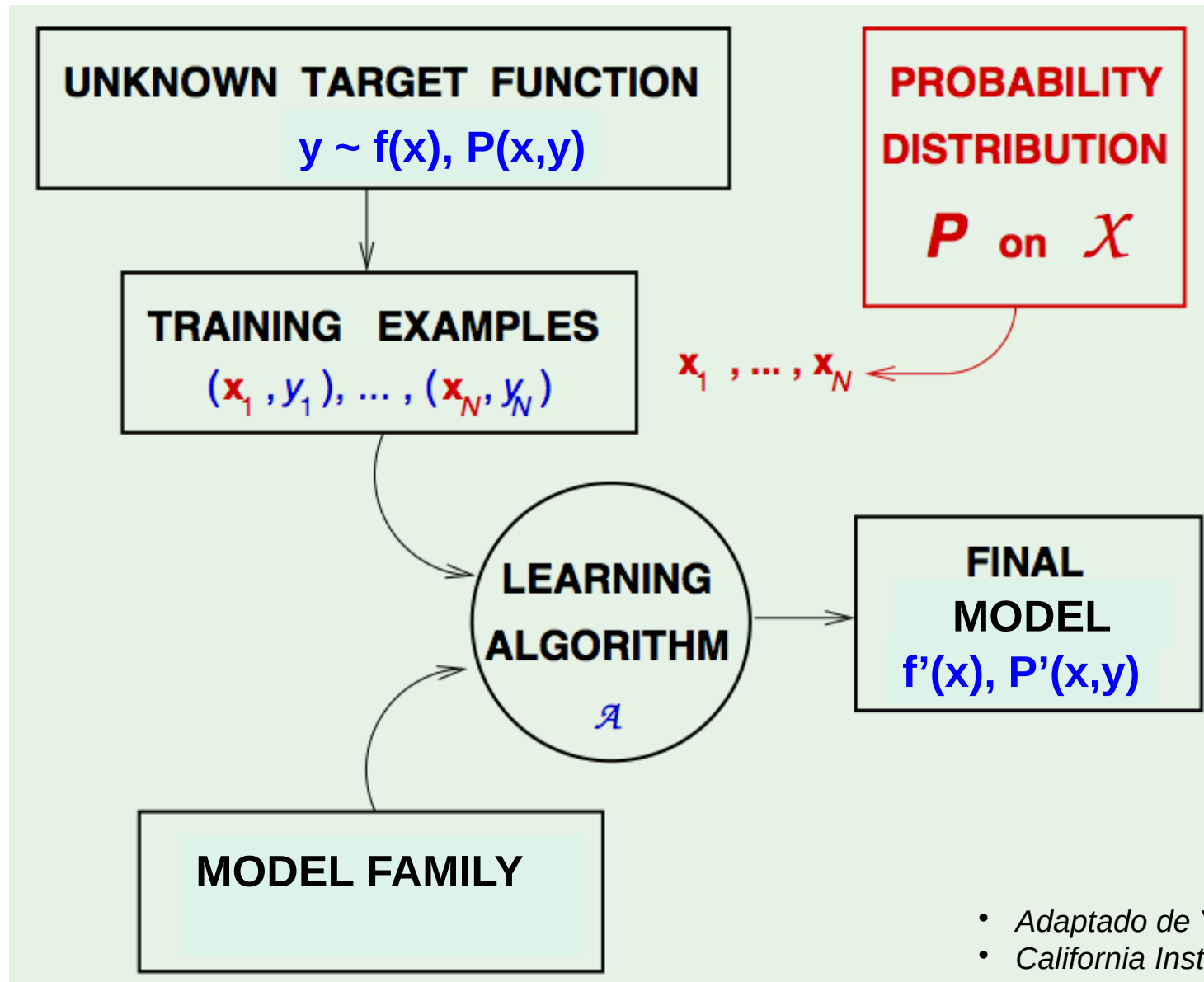
REGRESSION

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

Cross-Validation

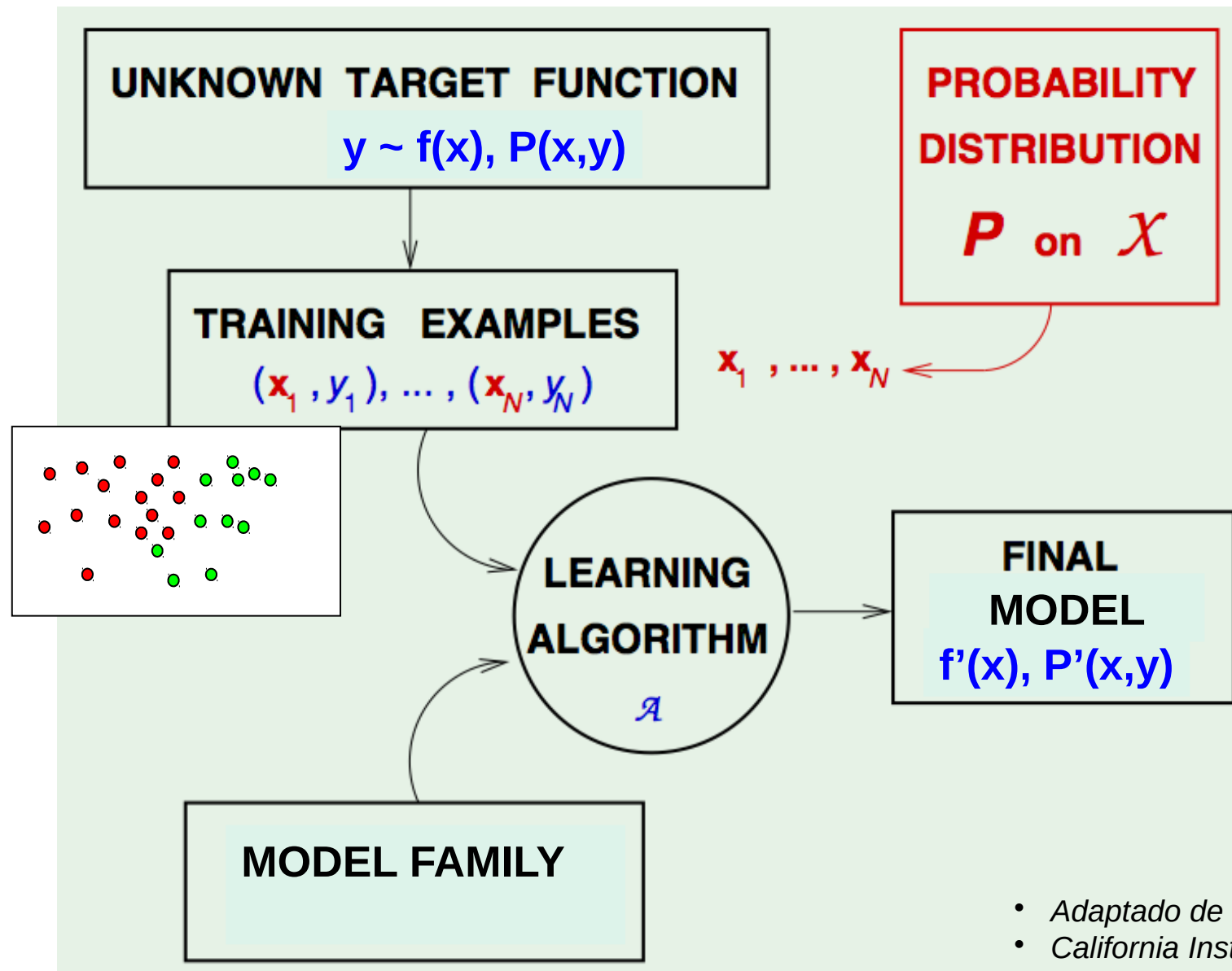
Data Mining: Data Modeling

Learning is the automatic process of building (adjusting) a model from a data set which is representative from the full population.



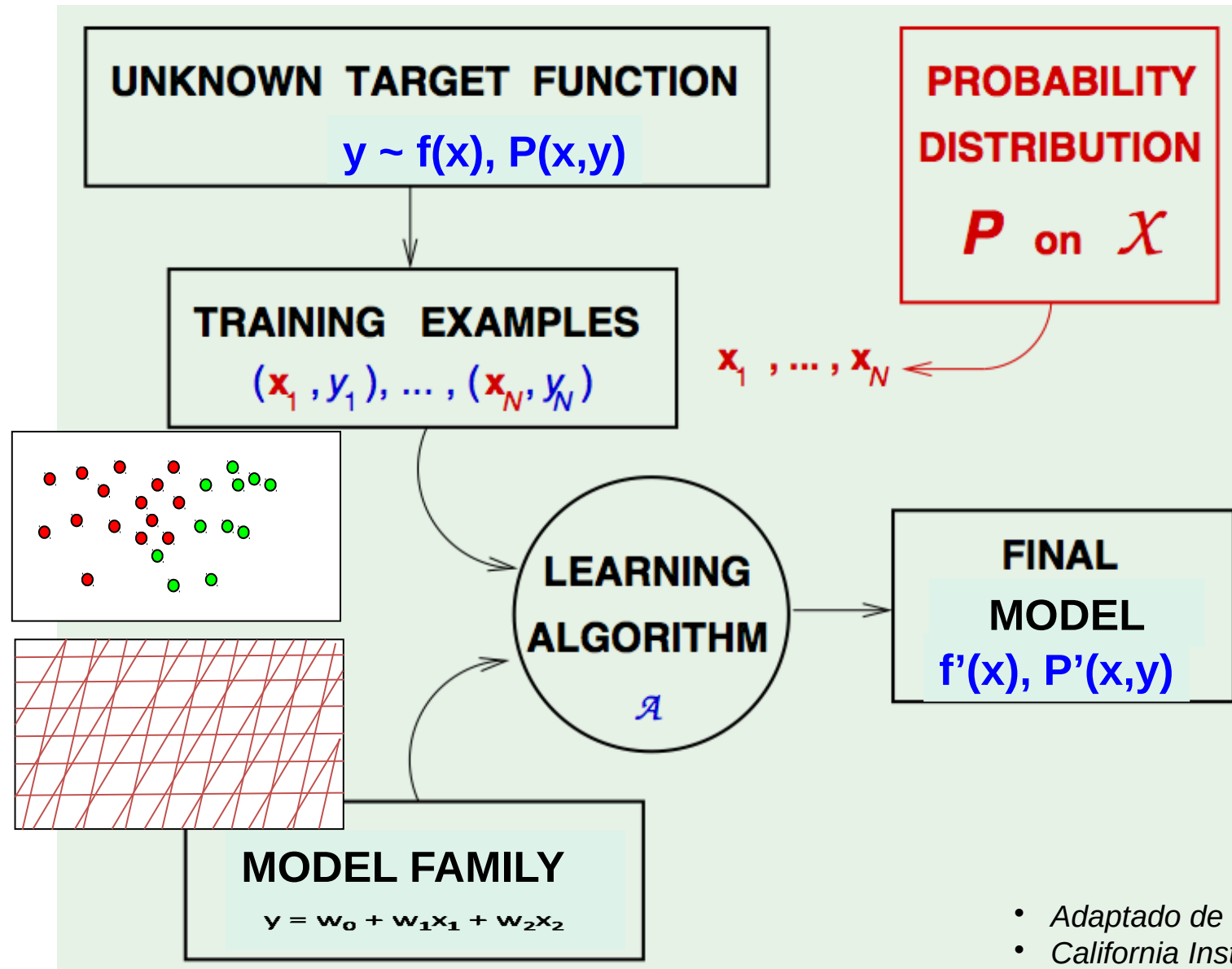
- Adaptado de Y.S. Abu-Mostafa
- California Institute of Technology

Learning is the automatic process of building (adjusting) a model from a data set which is representative from the full population.



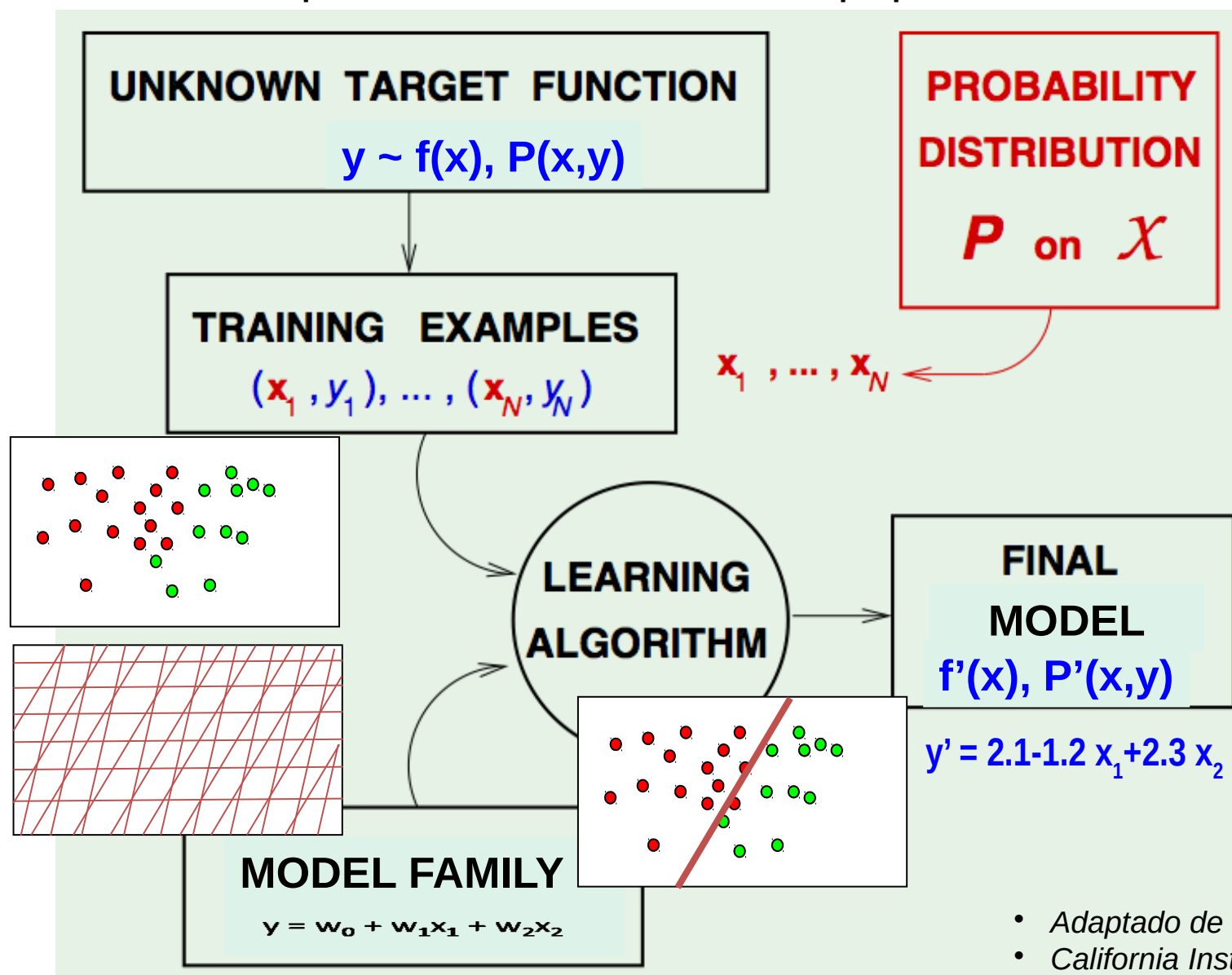
- Adaptado de Y.S. Abu-Mostafa
- California Institute of Technology

Learning is the automatic process of building (adjusting) a model from a data set which is representative from the full population.



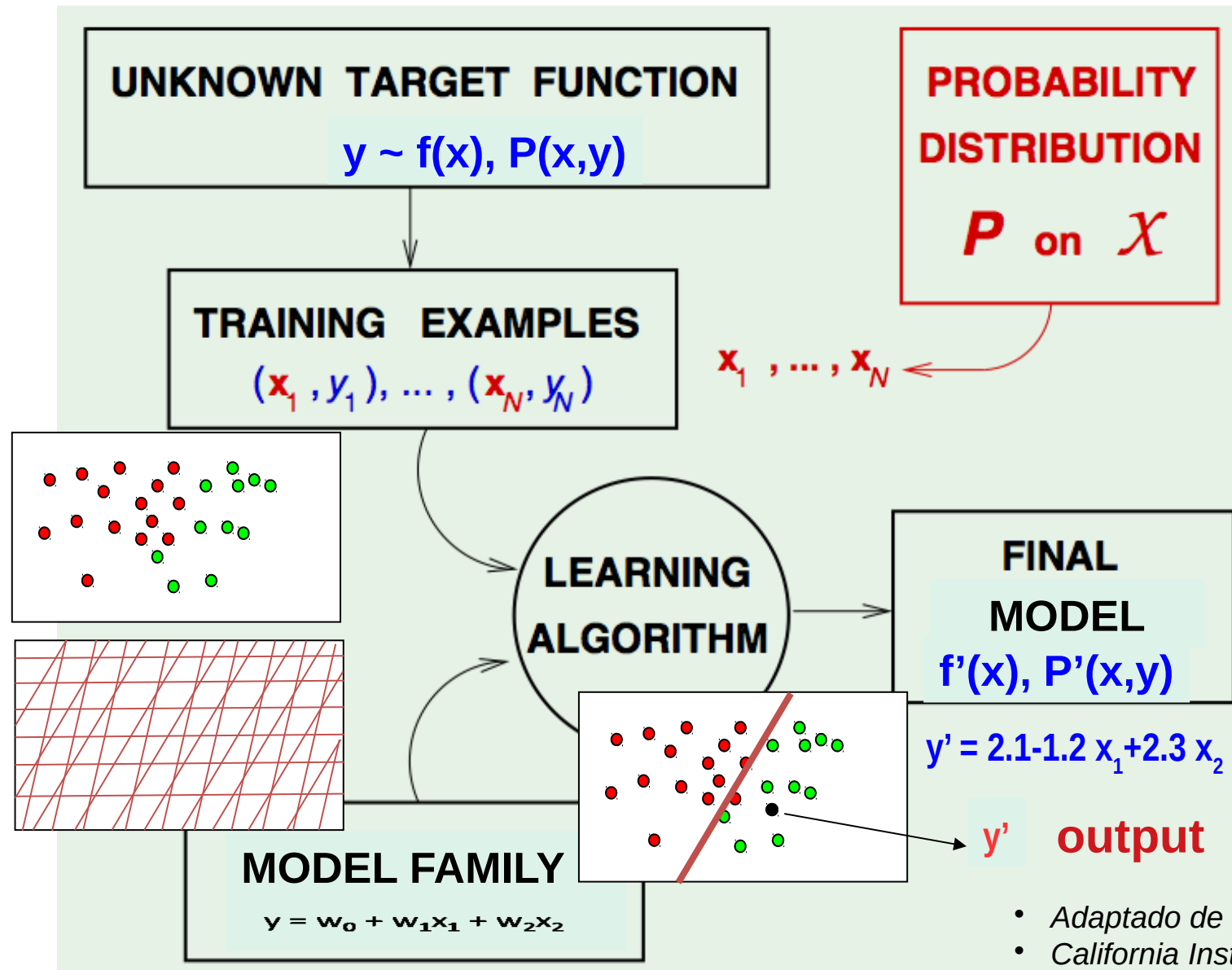
- Adaptado de Y.S. Abu-Mostafa
- California Institute of Technology

Learning is the automatic process of building (adjusting) a model from a data set which is representative from the full population.



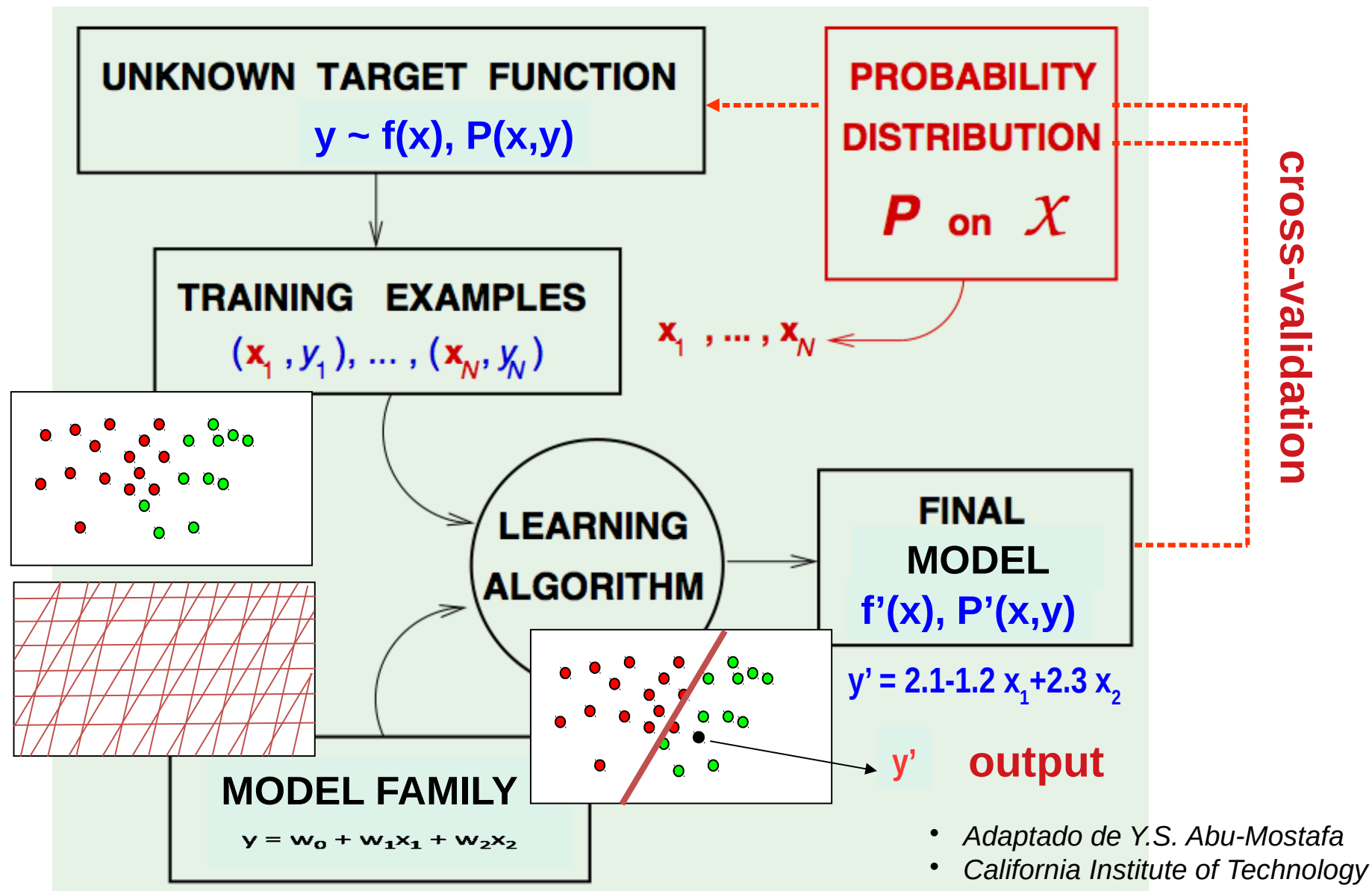
- Adaptado de Y.S. Abu-Mostafa
- California Institute of Technology

Learning is the automatic process of building (adjusting) a model from a data set which is representative from the full population.



- Adaptado de Y.S. Abu-Mostafa
- California Institute of Technology

Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).



UNKNOWN TARGET FUNCTION

$$y \sim f(x), P(x,y)$$

Finite/Incomplete sample

TRAINING EXAMPLES

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$$

$$\mathbf{x}_1, \dots, \mathbf{x}_N$$

PROBABILITY
DISTRIBUTION

P on \mathcal{X}

LEARNING
ALGORITHM

\mathcal{A}

FINAL
MODEL

$$f'(x), P'(x,y)$$

$$y' \quad P'(y|x)$$

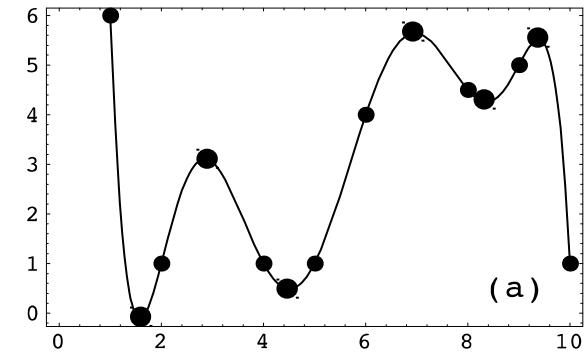
OUTPUT

MODEL FAMILY

\mathcal{H}

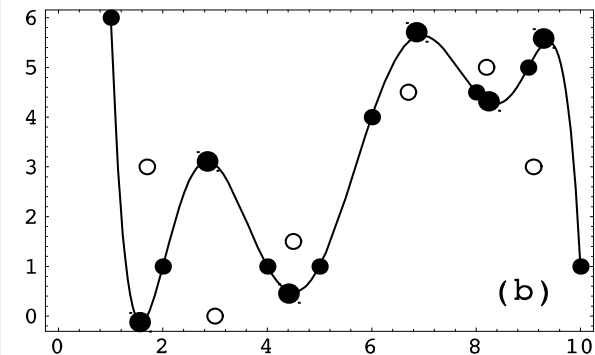
Adaptado de Y.S. Abu-Mostafa
California Institute of Technology

$$y = \sum_{i=0}^N (a_i * x^i)$$



New sample

Overfitting



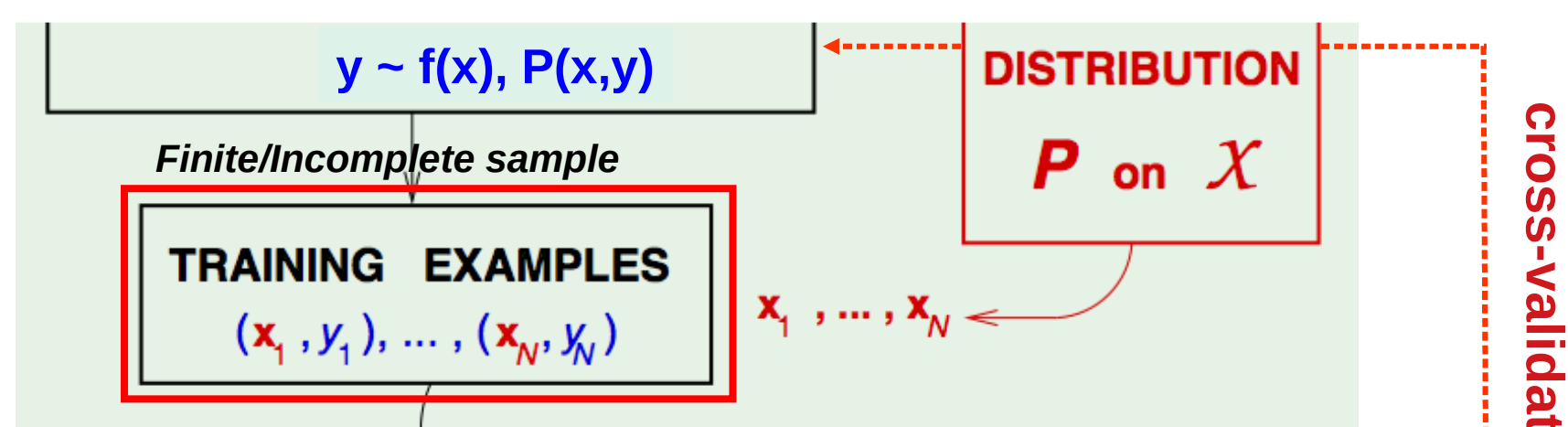
Trade-off between
bias and variance

Model complexity \leftrightarrow degrees of freedom

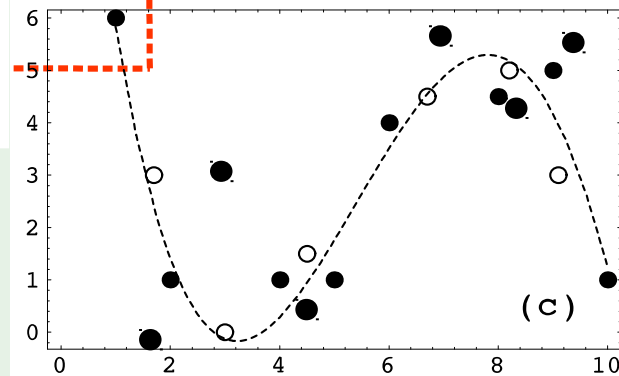
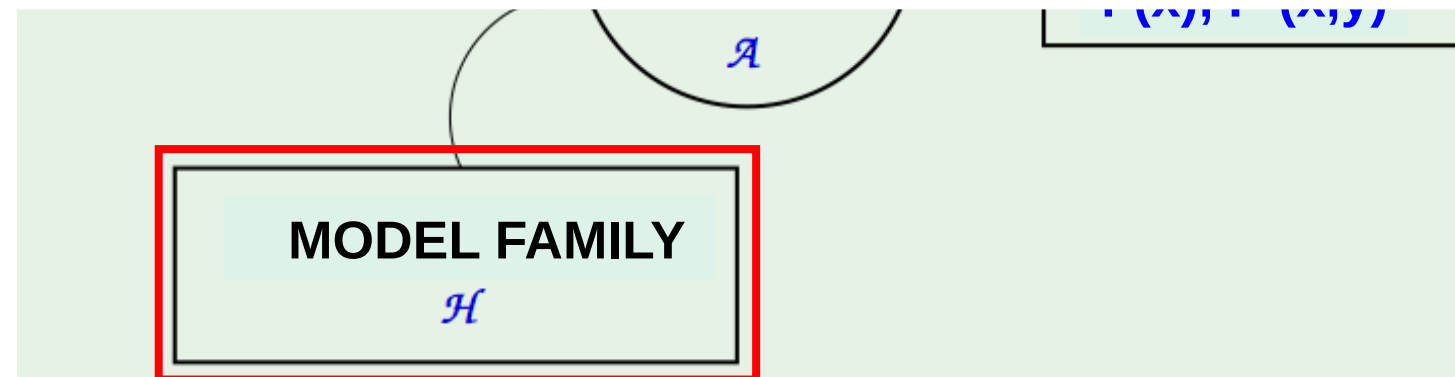
Cross-
Validation

LEARNING FROM DATA

Generalization is the most important feature for data driven systems: They must perform “well” when applied to new data (**cross-validation**).



Simplest models have better generalization properties and avoid the **overfitting** removing parameters/degree of freedom.



Trade-off between bias and variance

Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).

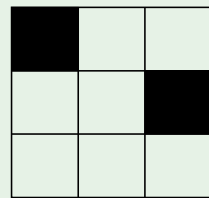
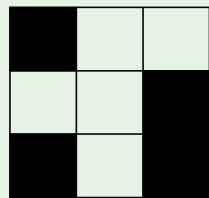
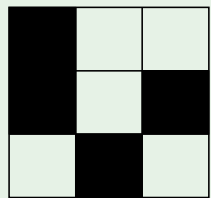
1. Can we make sure that $E_{out,1}(g)$ is close enough to $E_{in}(g)$?
2. Can we make $E_{in}(g)$ small enough?

Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).

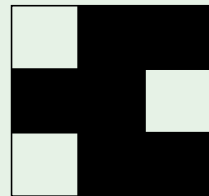
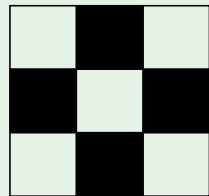
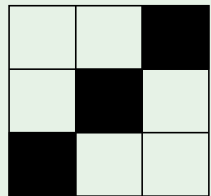
1. Can we make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$?
2. Can we make $E_{in}(g)$ small enough?

The (*in-sample*) error is the unique which can be estimated:

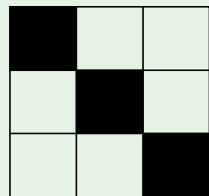
$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$



$f = -1$



$f = +1$



$f = ?$

		f	
		+ 1	- 1
h	+ 1	no error	false accept
	- 1	false reject	no error

$$E_{out}(h) = E(f, h)$$

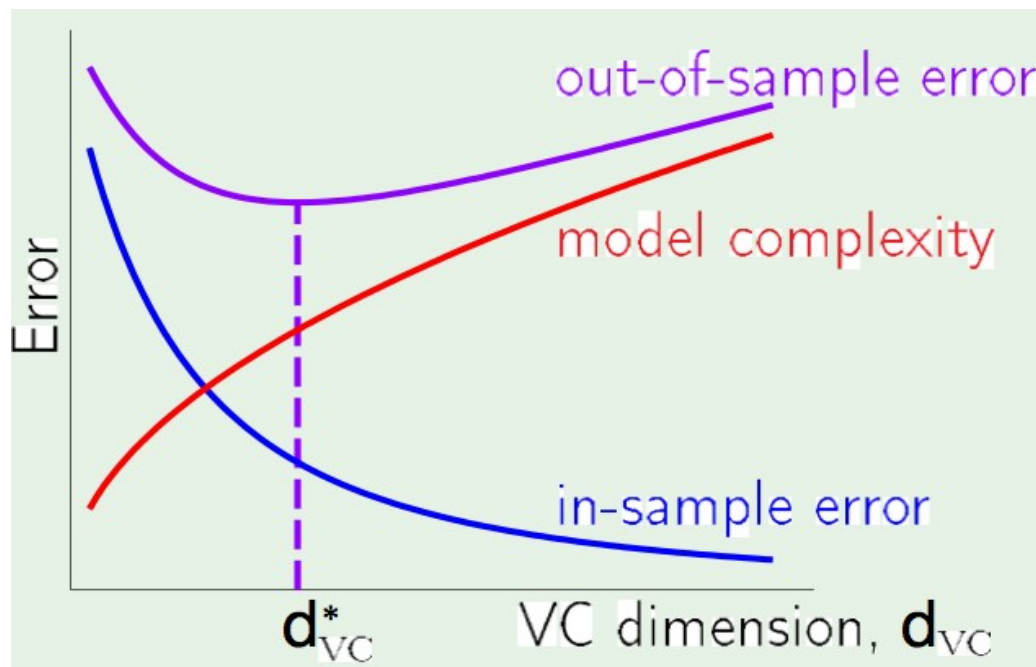
Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).

1. Can we make sure that $E_{out,1}(g)$ is close enough to $E_{in}(g)$?
2. Can we make $E_{in}(g)$ small enough?

The (*in-sample*) error is the unique which can be estimated:

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$

Vapnik-Chervonenkis (VC) Dimension



		f	
		+ 1	- 1
h	+ 1	no error	false accept
	- 1	false reject	no error

$$E_{out}(h) = E(f, h)$$

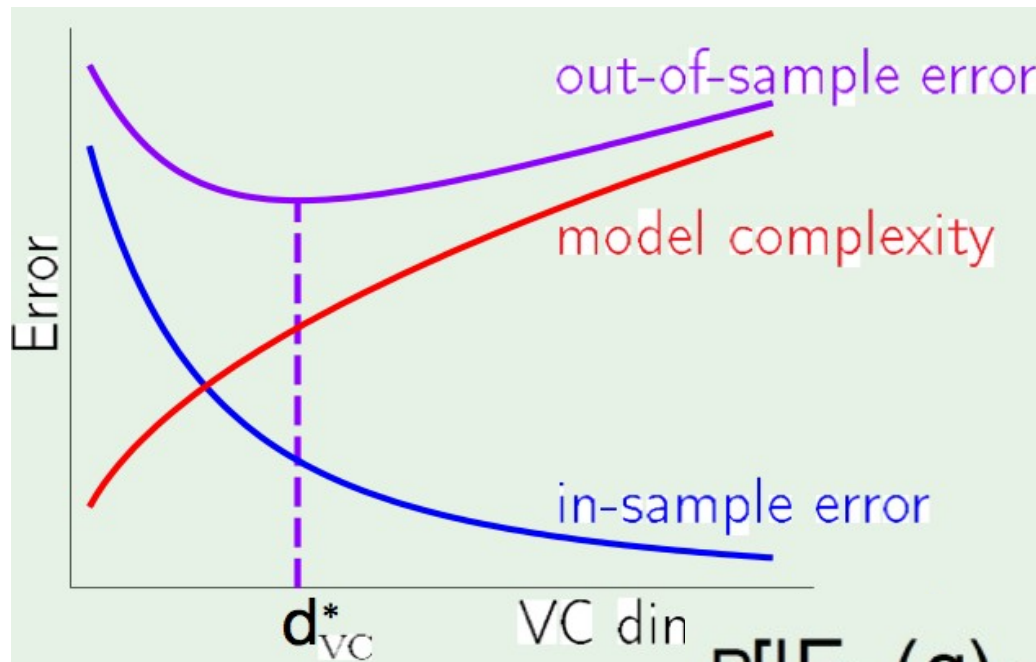
Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).

1. Can we make sure that $E_{out,1}(g)$ is close enough to $E_{in}(g)$?
2. Can we make $E_{in}(g)$ small enough?

The (*in-sample*) error is the unique which can be estimated:

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$

Vapnik-Chervonenkis (VC) Dimension



$$E_{out}(h) = E(f, h)$$

$$\mathbb{P}[|v - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

N =sample size

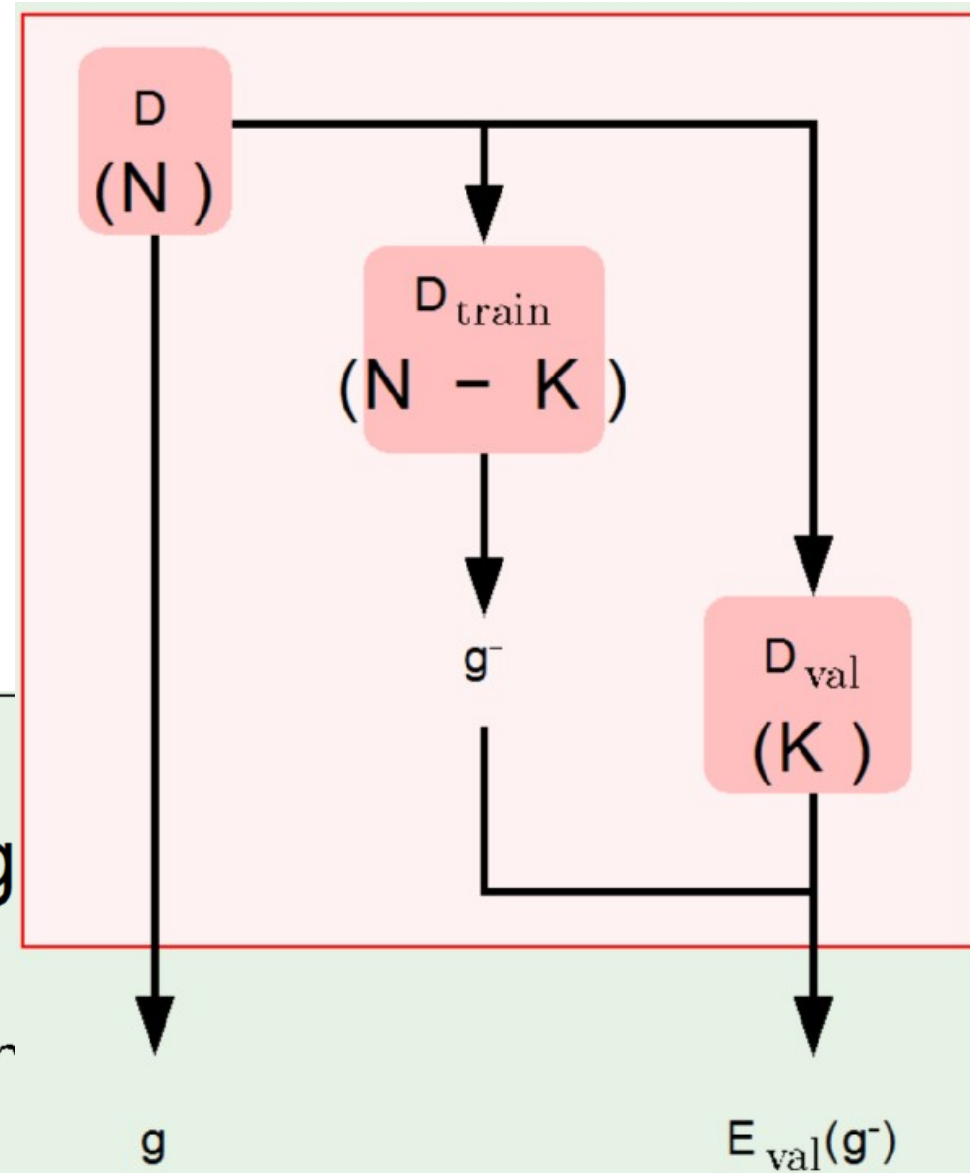
M =model complexity

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2M e^{-2\epsilon^2 N}$$

Generalization is the most important feature for data driven systems:
They must perform “well” when applied to new data (**cross-validation**).

The sample is divided in two subsets:
train and **test**.

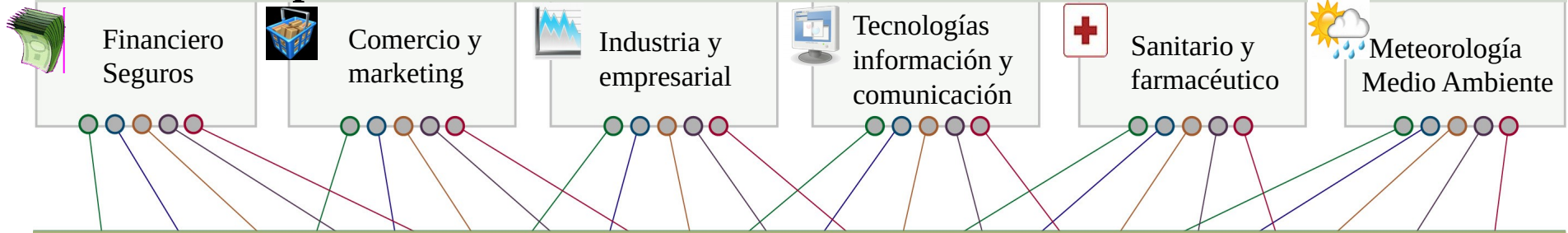
- Hold-out
- Leave-one-out
- K-fold



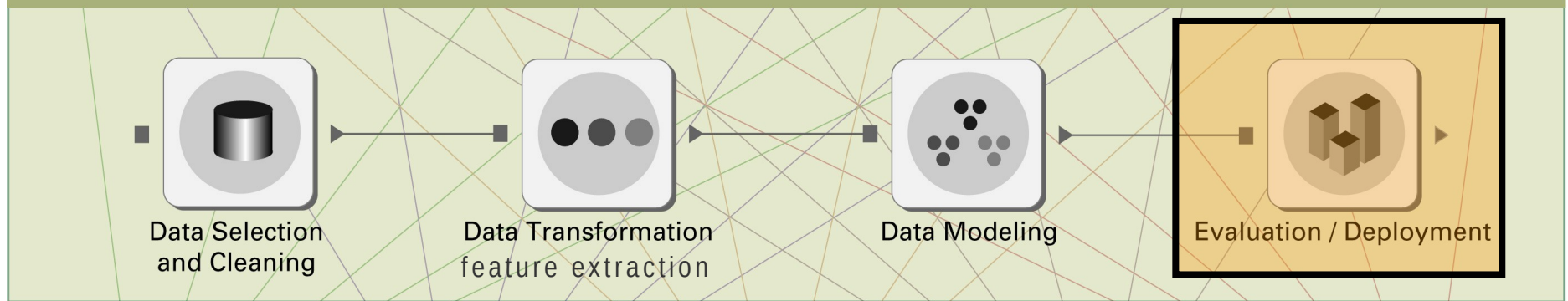
1. Can we make sure that $E_{\text{out}}(g_1)$

2. Can we make $E_{\text{in}}(g)$ smaller

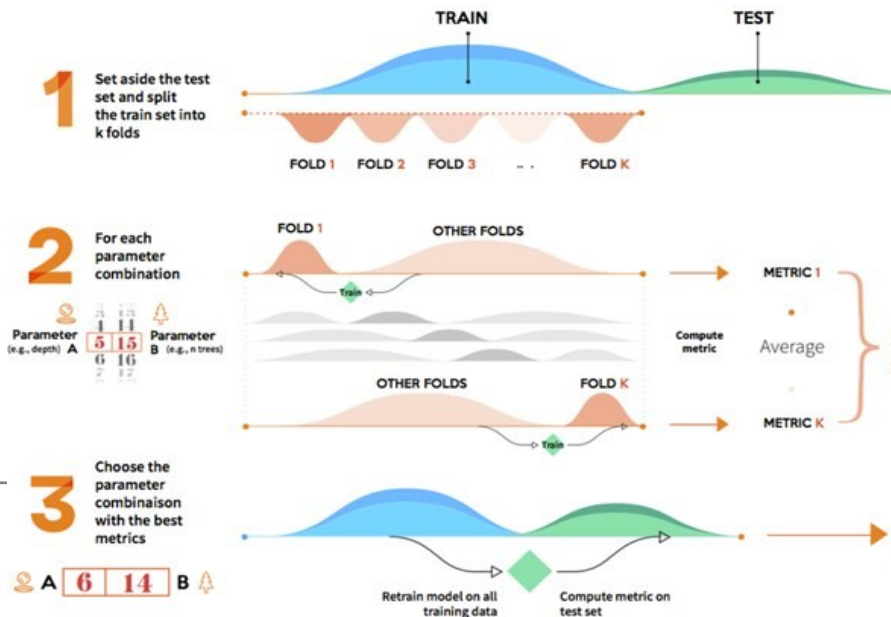
Sectores de aplicación



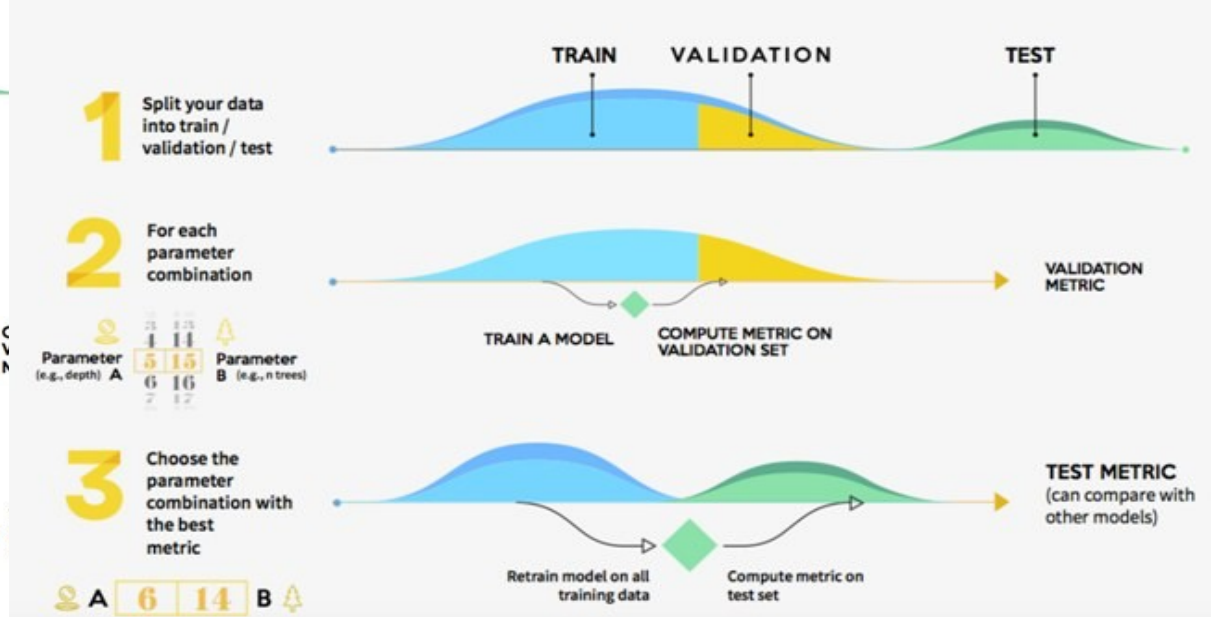
Proceso de Minería de Datos



K-FOLD STRATEGY



HOLDOUT STRATEGY



HOLDOUT STRATEGY

1 Split your data into train / validation

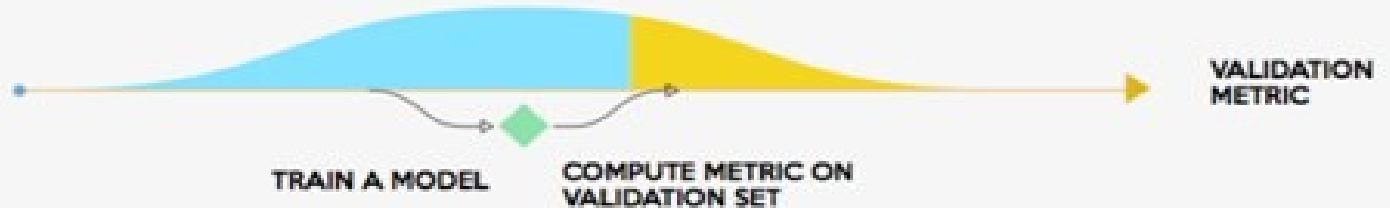


2 For each parameter combination

Parameter (e.g., depth) **A**

4	14
5	15
6	16
7	17

Parameter **B** (e.g., n trees)



Source: [Robert Kelley](#)

HOLDOUT STRATEGY

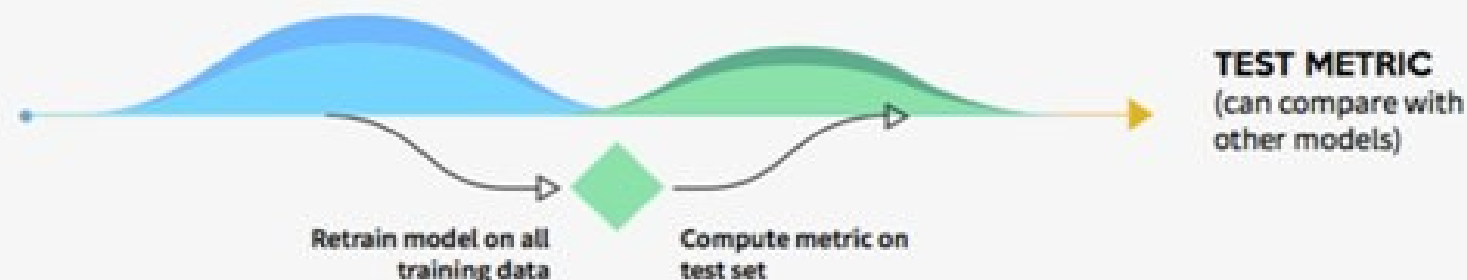
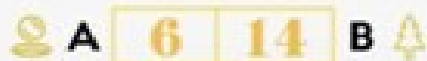
1 Split your data into train / validation / test



2 For each parameter combination



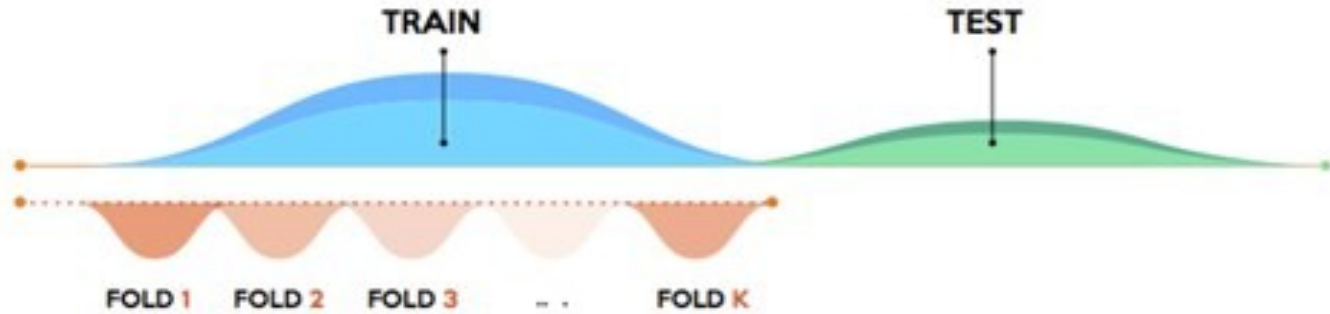
3 Choose the parameter combination with the best metric



Source: [Robert Kelley](#)

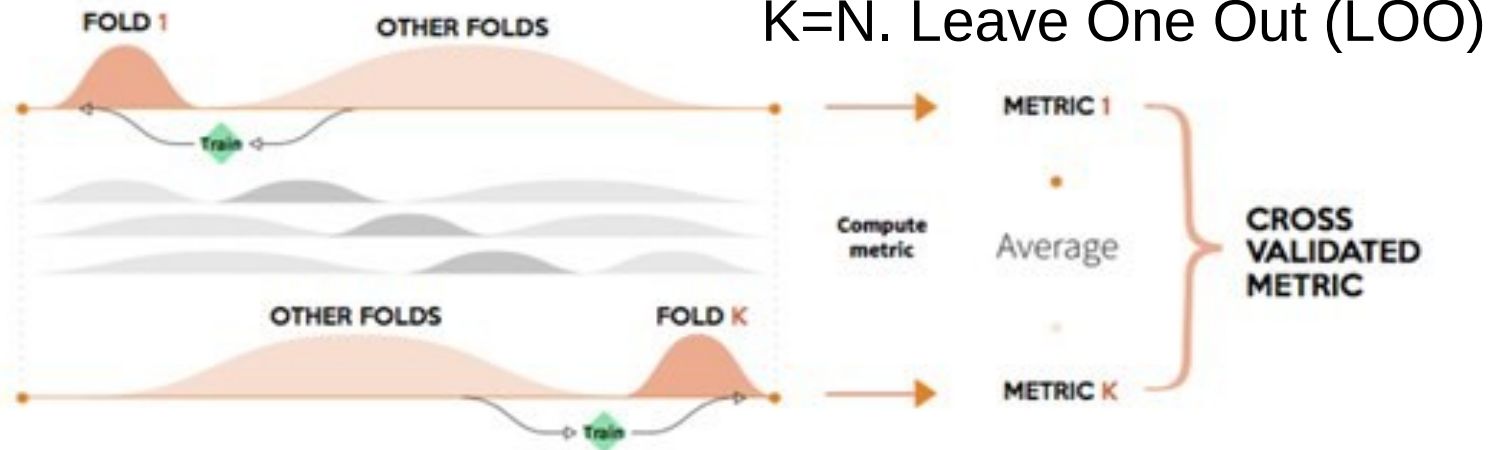
K-FOLD STRATEGY

1 Set aside the test set and split the train set into k folds



2 For each parameter combination

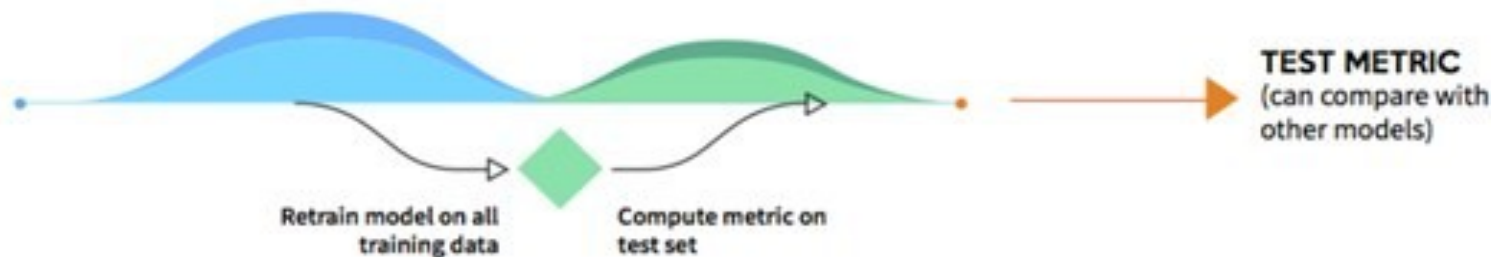
Parameter (e.g., depth)	A	1 2 3 4 5 6 7	11 12 13 14 15 16 17	Parameter (e.g., n trees)	B
		5	15		



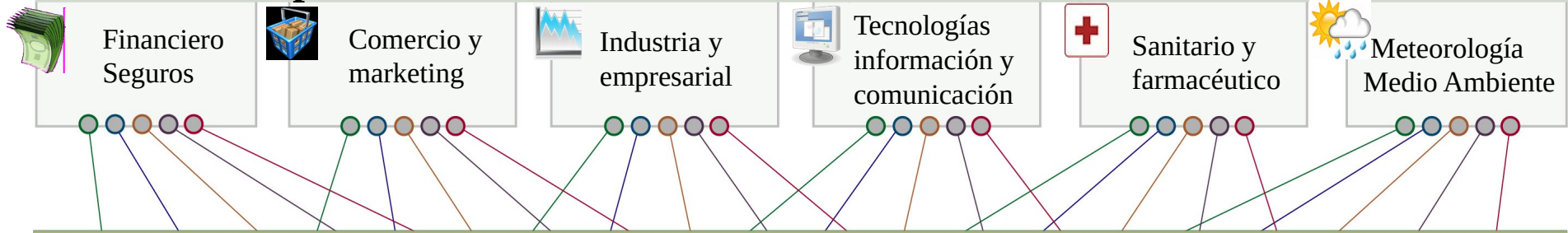
K=N. Leave One Out (LOO)

3 Choose the parameter combination with the best metrics

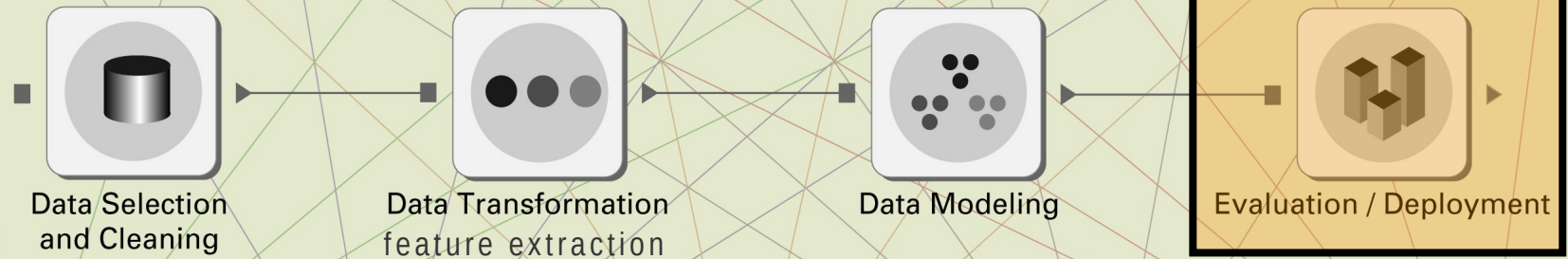
Parameter (e.g., depth)	A	6	14	Parameter (e.g., n trees)	B
----------------------------	---	---	----	------------------------------	---

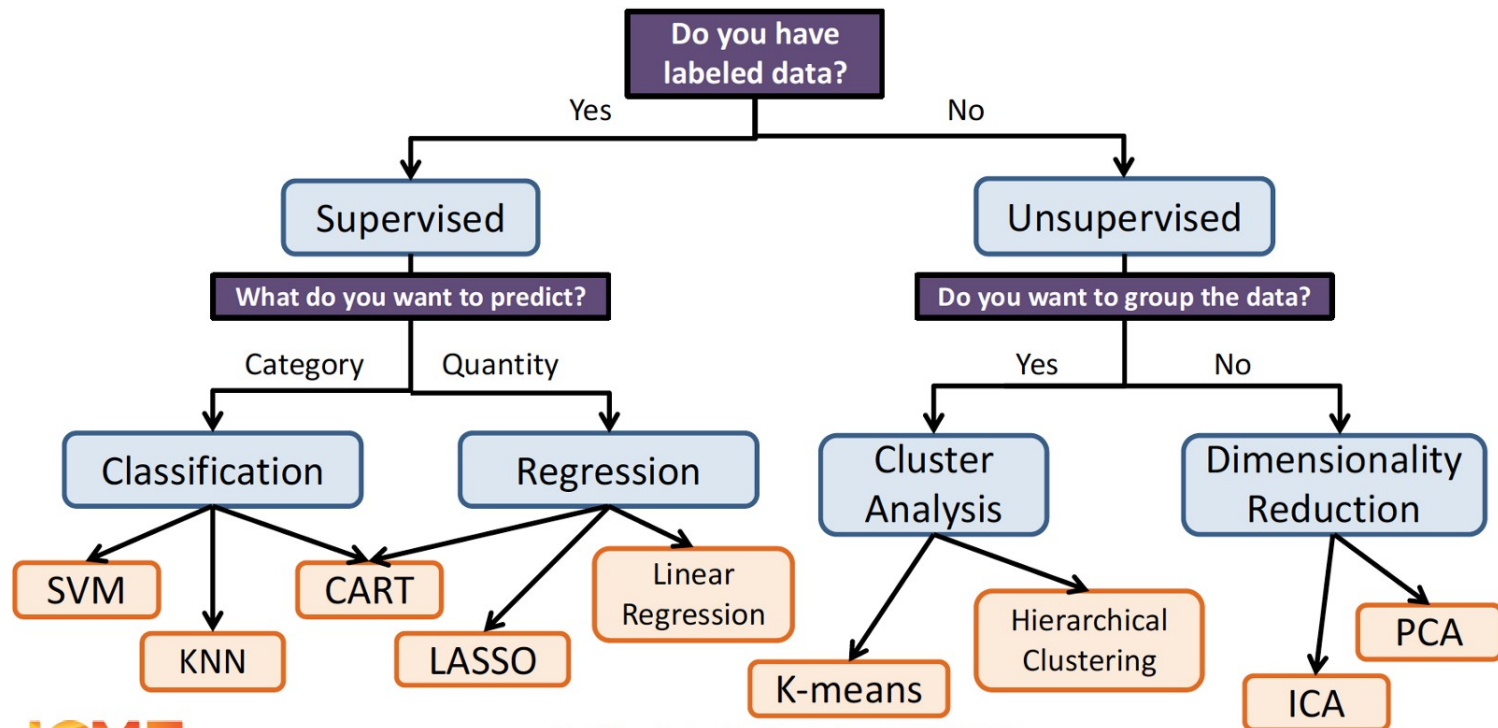
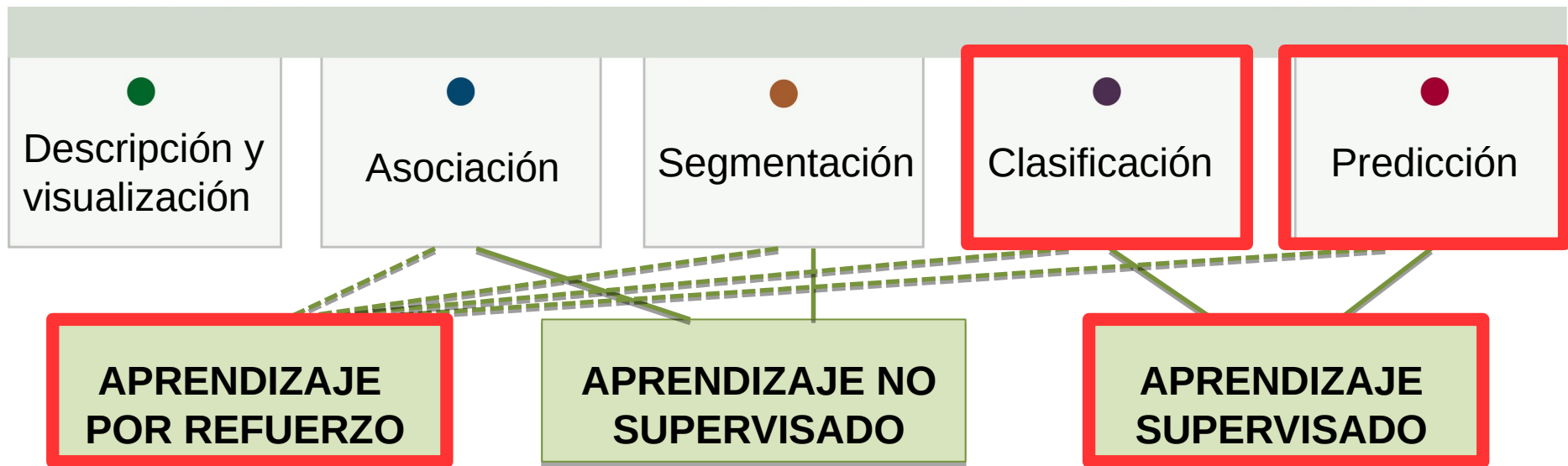


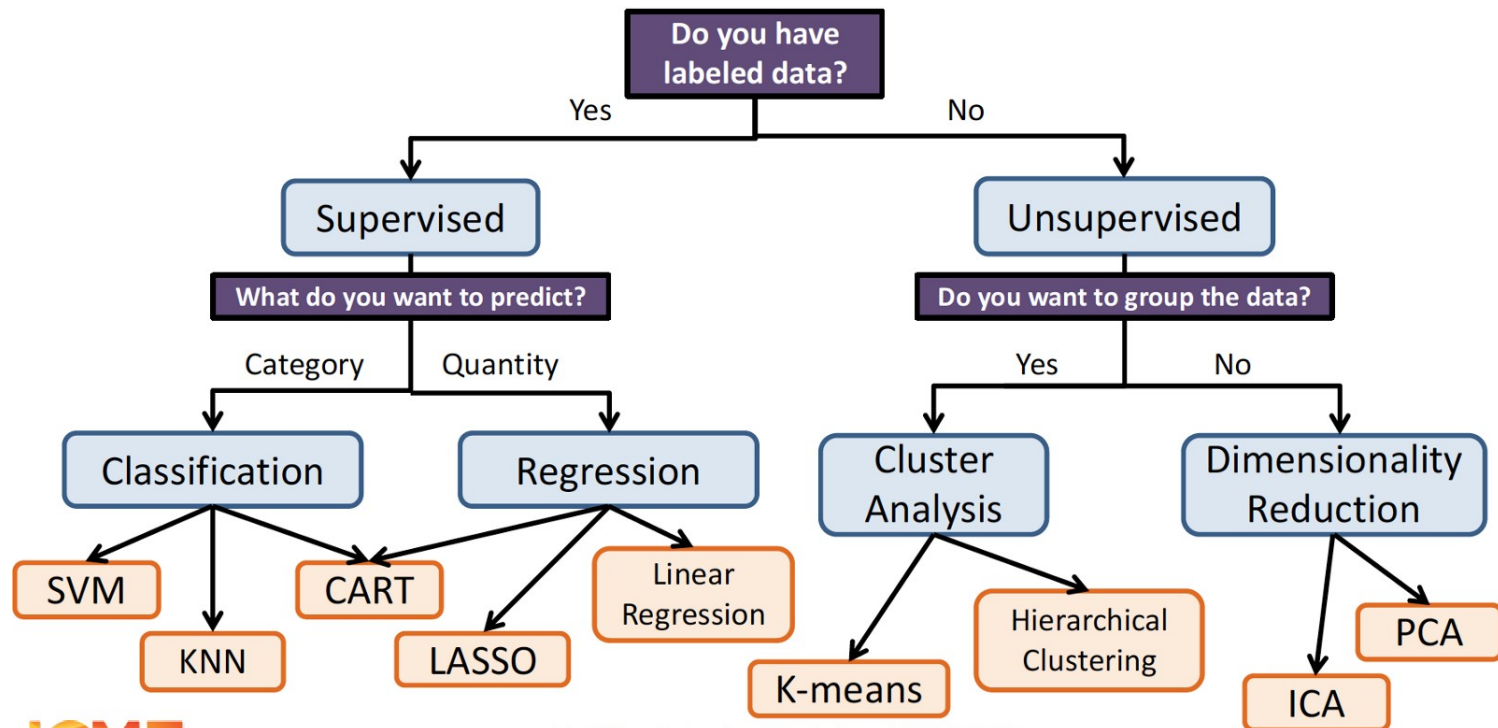
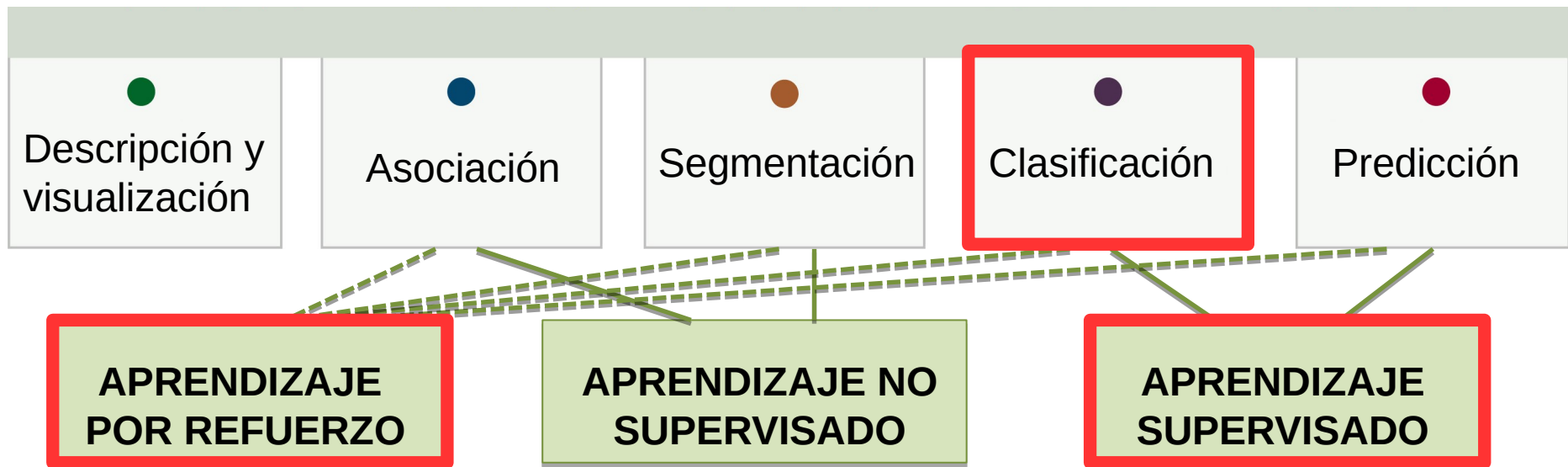
Sectores de aplicación



Proceso de Minería de Datos







		<u>True class</u>	
		p	n
<u>Predicted class</u>	Y	True Positives	False Positives
	N	False Negatives	True Negatives
Column totals:		P	N

False Alarm Rate (**FAR**) Hit Rate (**HIR**)

$$\text{fp rate} = \frac{FP}{N}$$

$$\text{tp rate} = \frac{TP}{P}$$

ACCURACY FUNCTION	$(TP + TN) / (P + N)$
PRECISION FUNCTION	$TP / (TP + FP)$
SPECIFICITY FUNCTION	$TN / (FP + TN)$
SENSITIVITY FUNCTION	$TP / (TP + FN)$

FAR = 1-specificity

HIR = sensitivity

Fawcett, T. (2006) An introduction to ROC analysis, In Pattern Recognition Letters, 27, 861-874, <https://doi.org/10.1016/j.patrec.2005.10.010>.

		<u>True class</u>	
		p	n
<u>Predicted class</u>	Y	True Positives	False Positives
	N	False Negatives	True Negatives

Column totals:

P

N

Which systems yield?

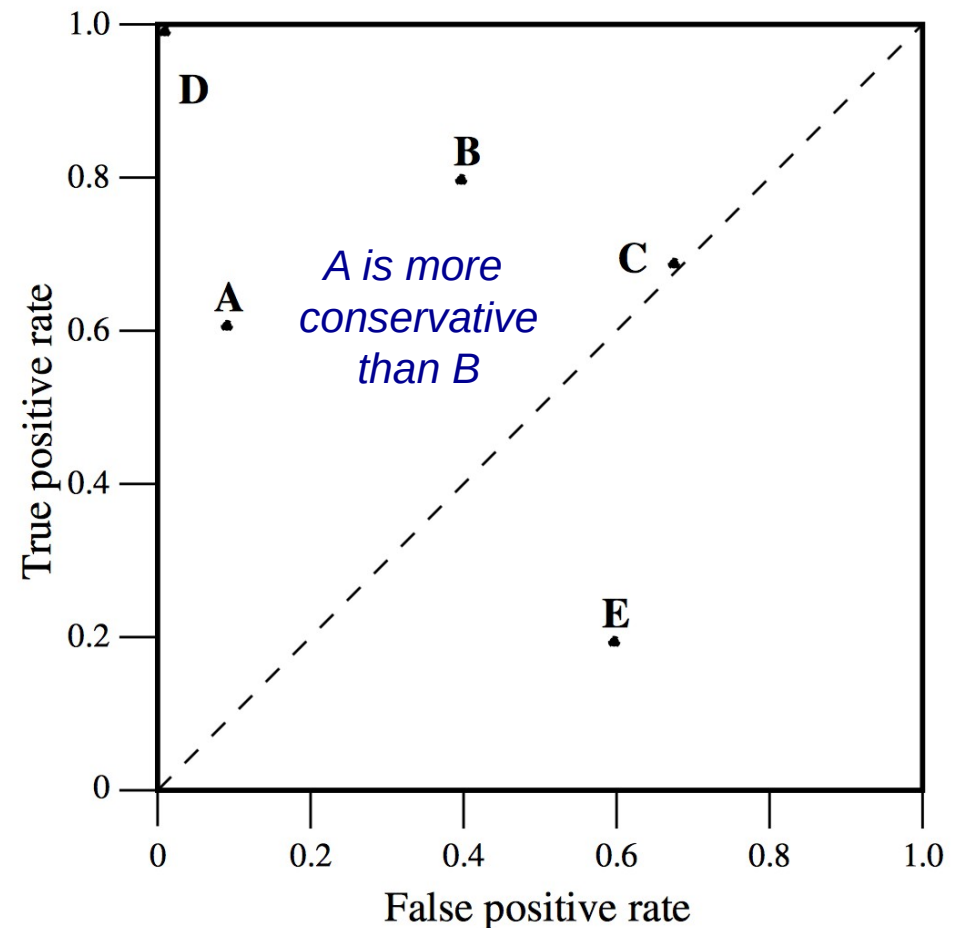
HIR = FAR = 0 → Never predicting

HIR = FAR = 1 → Always predicting

False Alarm Rate (**FAR**) Hit Rate (**HIR**)

$$\text{fp rate} = \frac{FP}{N}$$

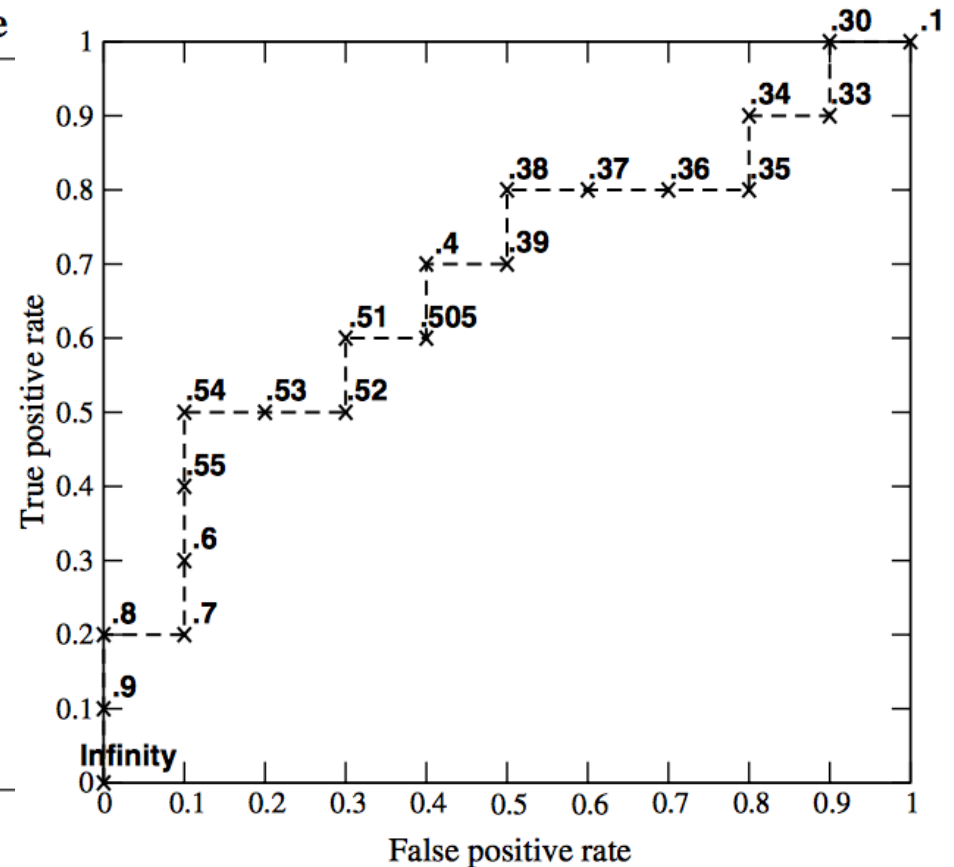
$$\text{tp rate} = \frac{TP}{P}$$



Summarizes the performance of the system over all possible probability thresholds.

```
library(pROC)
obs<-c(rep(0,50),rep(1,50));
prd<-obs+2*(runif(100)-0.5);
prd[which(prd<0)]<-0; prd[which(prd>1)]<-1;
plot(roc(obs,prd), print.auc=TRUE)
hist(prd)
```

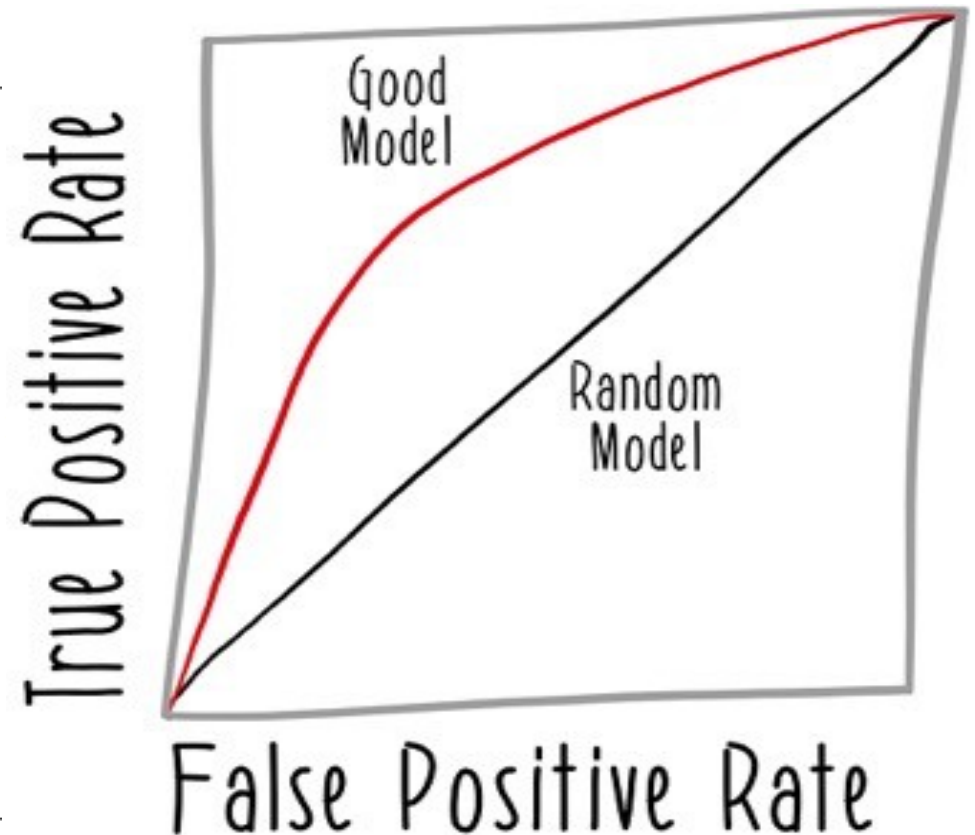
Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1



Summarizes the performance of the system over all possible probability thresholds.

```
library(pROC)
obs<-c(rep(0,50),rep(1,50));
prd<-obs+2*(runif(100)-0.5);
prd[which(prd<0)]<-0; prd[which(prd>1)]<-1;
plot(roc(obs,prd), print.auc=TRUE)
hist(prd)
```

Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1



<https://www.kdnuggets.com/2018/01/machine-learning-model-metrics.html>

- **Accuracy:** Overall, how often is the classifier correct?
 - $(TP+TN)/total = (100+50)/165 = 0.91$
- **Misclassification Rate:** Overall, how often is it wrong?
 - $(FP+FN)/total = (10+5)/165 = 0.09$
 - equivalent to 1 minus Accuracy
 - also known as "Error Rate"
- **True Positive Rate:** When it's actually yes, how often does it predict yes?
 - $TP/actual\ yes = 100/105 = 0.95$
 - also known as "Sensitivity" or "Recall"
- **False Positive Rate:** When it's actually no, how often does it predict yes?
 - $FP/actual\ no = 10/60 = 0.17$
- **Specificity:** When it's actually no, how often does it predict no?
 - $TN/actual\ no = 50/60 = 0.83$
 - equivalent to 1 minus False Positive Rate
- **Precision:** When it predicts yes, how often is it correct?
 - $TP/predicted\ yes = 100/110 = 0.91$
- **Prevalence:** How often does the yes condition actually occur in our sample?
 - $actual\ yes/total = 105/165 = 0.64$

HIR
(Hit rate)

FAR
(False alarm rate)

For which systems
do the following
equalities hold?:

$HIR = FAR = 0$

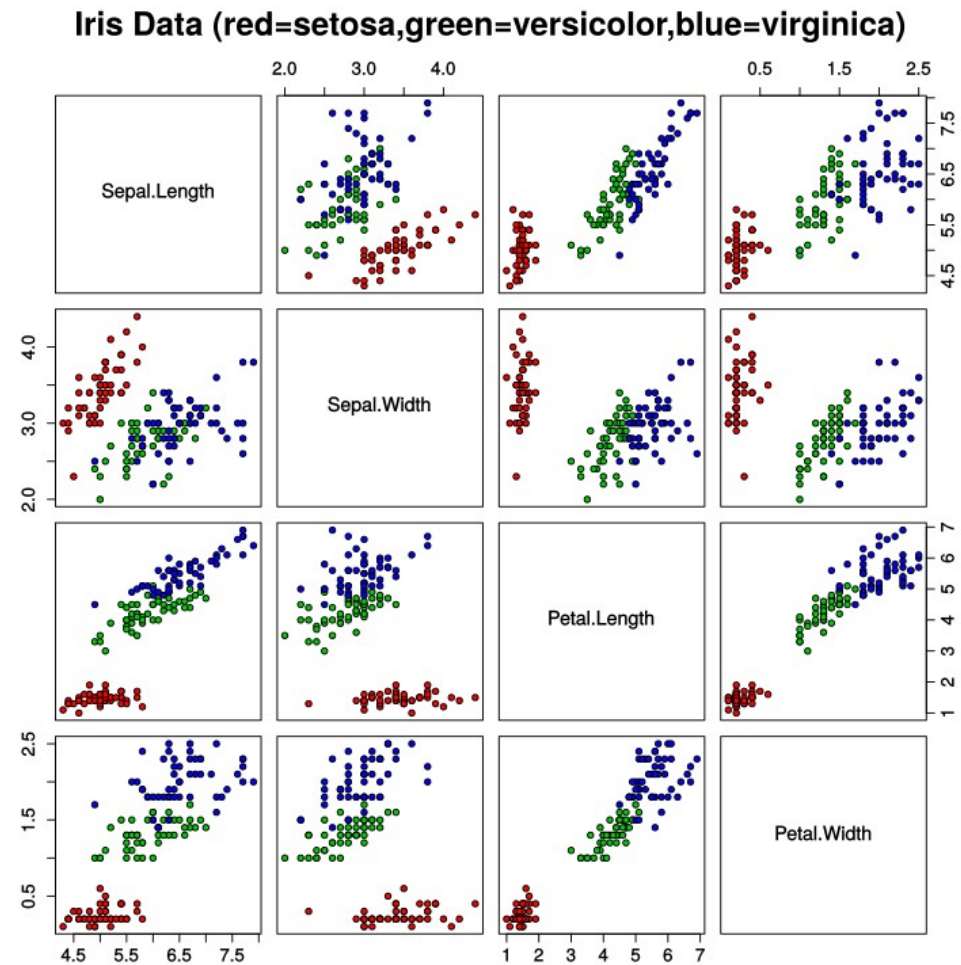
$HIR = FAR = 1$

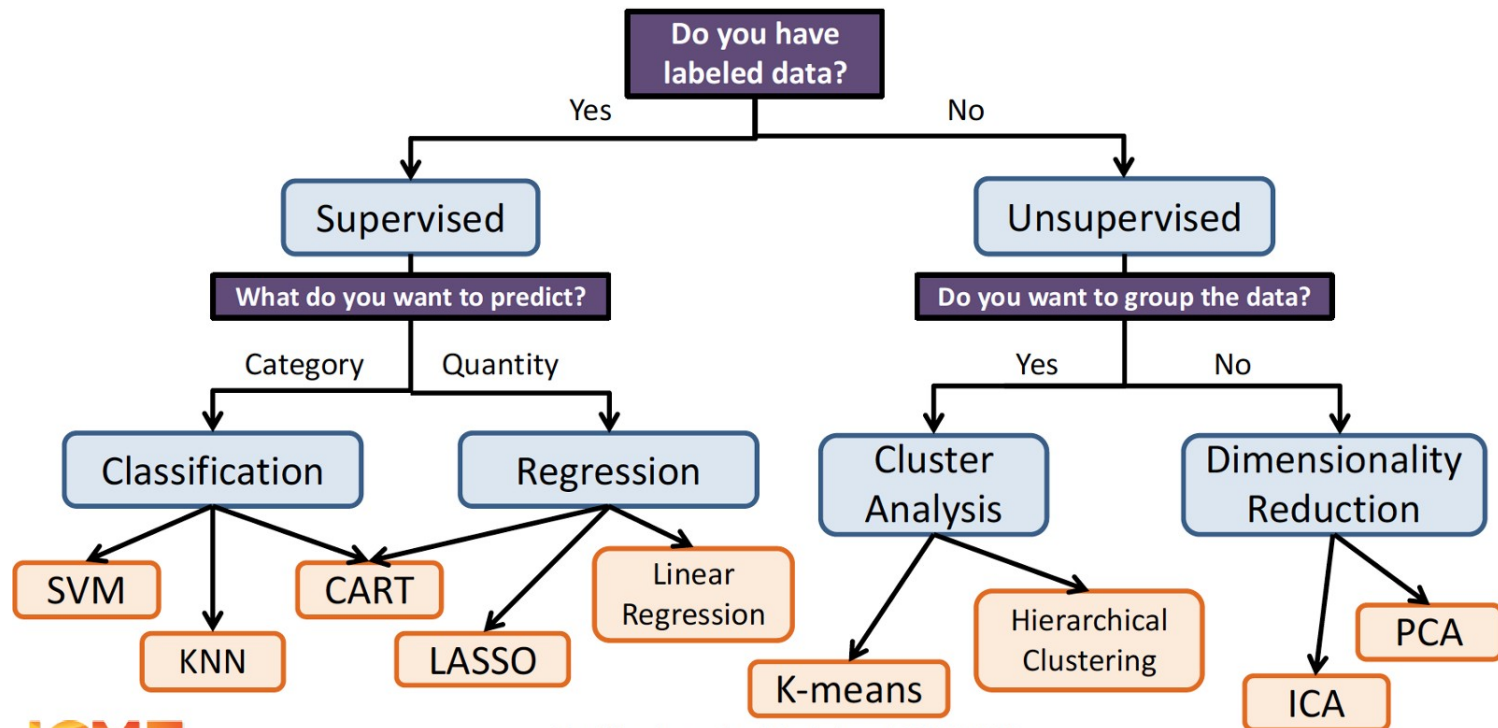
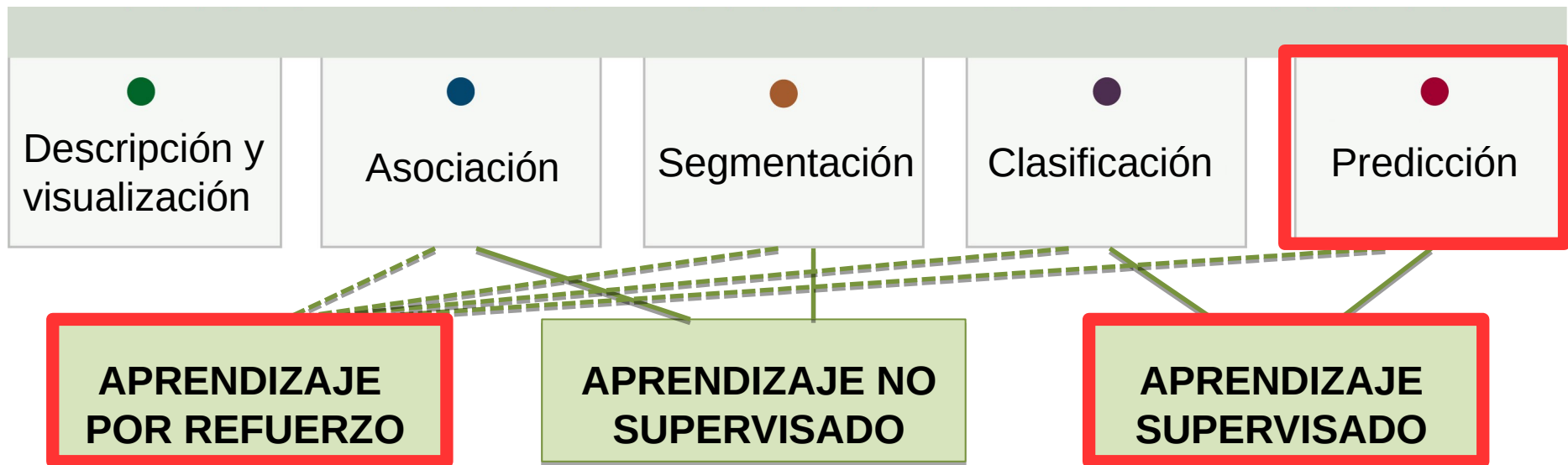
```

data(iris)
fitControl <- trainControl(method="none",
                           number=1,
                           repeats=1,
                           verboseIter=TRUE)

modelFit <- train(Species ~ ., data=iris, method="knn", trControl=fitControl)
pred <- predict(modelFit, newdata = iris[,-5])
acc<-confusionMatrix(iris$Species,pred)
print(acc)

```





Model accuracy (training and validation).

Some models are trained using an **empirical error (cost) function**, which measures **model accuracy** as the difference between the predicted and the actual value. In this case, this is a natural **validation measure**.

- This cost function could be anything:
 - Sum of absolute errors: $J = \sum |y - u|$.
 - Sum of square errors: $J = \sum (y - u)^2$.
 - As long as the minimum occurs when the distributions are the same, in theory it would work.
- One good idea is that u represents the parameters of the distribution of y .
 - Rationale: often natural processes are fuzzy, and any input might have a range of outputs.
 - This approach also gives a smooth measure of how accurate we are.
 - The maximum likelihood principle says that: $\theta_{\text{ML}} = \arg \max_{\theta} p(y; u)$
 - Thus we want to minimize: $J = -p(y; u)$
 - For i samples: $J = -\prod_i p(y_i; u)$
 - Taking log both sides: $J' = -\sum_i \log p(y_i; u)$.
 - This is called cross-entropy.
- Applying the idea for: $y \sim \text{Gaussian}(\text{center} = u)$:
 - $p(y; u) = e^{-(y-u)^2}$.
 - $J = -\sum \log e^{-(y-u)^2} = \sum (y - u)^2$
 - This motivates sum of squares as a good choice.

Correlation (Pearson, Spearman)

Model performance: Validation diagnostics and metrics.

There are several domain-dependent diagnostics (computed separately for prediction 'p' and observation 'o') and metrics/errors for validating model performance.

Distributional consistency: evaluates the model capability to reproduce the distribution of the observed data.

- **Bias** = $\text{mean } p - \text{mean } o$
- **Variance ratio** = $\text{var } p / \text{var } o$
- **Distributional similarity:** ks-score, Von Misses, pdf-score, etc.

The **quantile-quantile plot** is a typical tool to evaluate, in a graphical way, the distributional similarity of the order statistics (e.g. **percentiles**).

Different diagnostics for different fields.

Accuracy: assess the correspondence of the simulated and observed sequences. Two typical scores are usually used: Root Mean Square Error (**RMSE**) and the (Pearson/Spearman/Kendall) **Correlation**.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Distributional consistency: evaluates the model capability to reproduce the distribution of the observed data. The most popular are the **bias** (mean difference) or the **ratio of variances/standard deviation**. In addition, there are hypothesis tests to evaluate in a global way the similarity of the observed and simulated series (e.g. **Kolmogorov-Smirnov**, **Perkins**, **Von Mises**, etc).

The **quantile-quantile plot** is a typical tool to evaluate, in a graphical way, the distributional similarity of the order statistics (e.g. **percentiles**).

Example with R:

?qqplot

require(graphics)

y<-rt(200,df=5)

qqnorm(y)

qqline(y,col=2)

qqplot(y,rt(300,df=5))

Accuracy: assess the correspondence of the simulated and observed sequences. Two typical scores are usually used: Root Mean Square Error (**RMSE**) and the (Pearson/Spearman/Kendall) **Correlation**.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Distributional consistency: evaluates the model capability to reproduce the distribution of the observed data. The most popular are the **bias** (mean difference) or the **ratio of variances/standard deviation**. In addition, there are hypothesis tests to evaluate in a global way the similarity of the observed and simulated series (e.g. **Kolmogorov-Smirnov**, **Perkins**, **Von Mises**, etc).

The **quantile-quantile plot** is a typical tool to evaluate, in a graphical way, the distributional similarity of the order statistics (e.g. **percentiles**).

Example with R:

?qqplot

require(graphics)

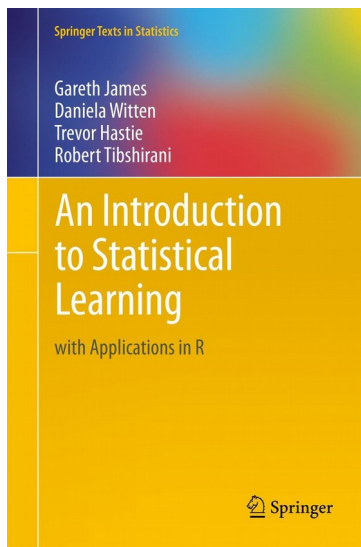
y<-rt(200, df=5)

qqnorm(y)

qqline(y, col=2)

qqplot(y, rt(300, df=5))

How to create and use our own functions, including validation measures in R?



An Introduction to Statistical Learning: With Applications in R

James, G., Witten, D., Hastie, T., Tibshirani, R.

Springer (2013)

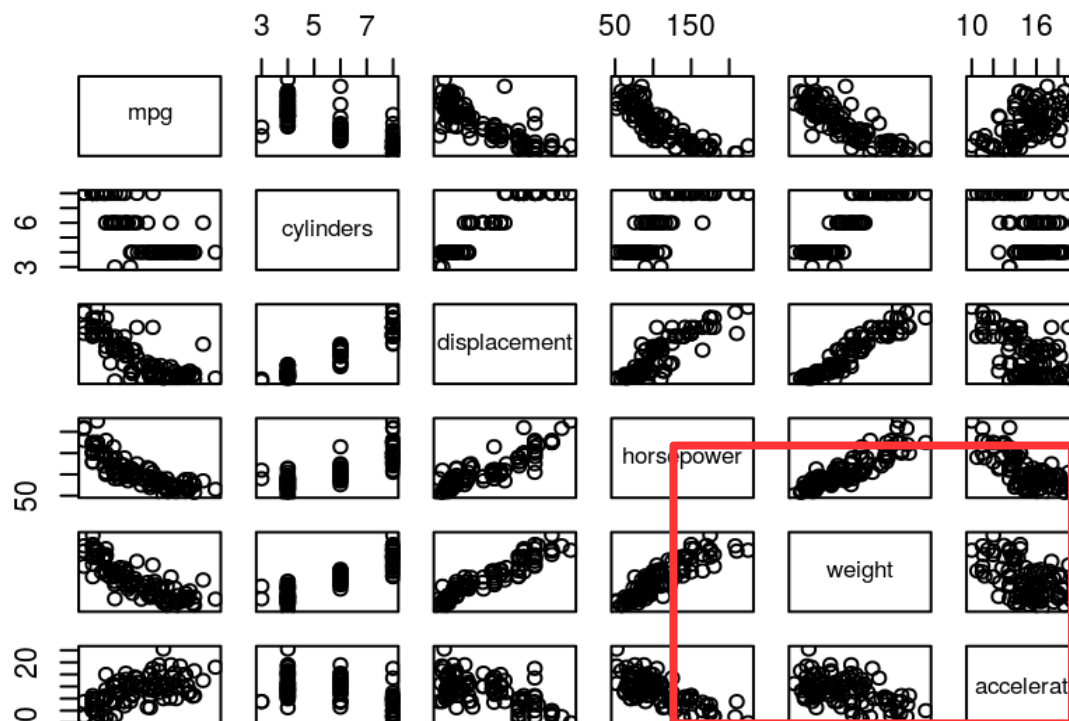
<http://www-bcf.usc.edu/~gareth/ISL>

```
install.packages("ISLR")
library("ISLR")
library(help = "ISLR")
```

```
> data(Auto)
> str(Auto)
```

```
'data.frame':    392 obs. of  9 variables:
 $ mpg          : num  18 15 18 16 17 ...
 $ cylinders    : num   8  8  8  8  8  8  8 ...
 $ displacement: num  307 350 318 304 ...
 $ horsepower   : num  130 165 150 150 ...
 $ weight       : num  3504 3693 3436 ...
 $ acceleration: num   12 11.5 11 12 ...
 $ year         : num   70  70  70  70  70 ...
 $ origin       : num    1  1  1  1  1  1  1 ...
 $ name         : Factor w/ 304 levels ...
```

```
> pairs(Auto)
```



CARET (C**l**Assification and R**E**gresion Training) is a wrapper of a number of standard machine learning packages which performs model tuning (optimization of the model parameters) and cross-validation strategies.
<http://topepo.github.io/caret/index.html>

```
> modelLookup(model = "lm")
  model parameter      label forReg forClass probModel
  lm intercept intercept   TRUE   FALSE   FALSE
```

```
trainControl(method , number, ...)
  method: "none", "cv", "LOOCV"
  number: For "cv" (2 => hold-out, 10 => 10-fold)
```

```
> ctrl <- trainControl(method = "LOOCV")
> mod <- train(weight ~ horsepower,
               data = Auto,
               method = "lm",
               trControl = ctrl)
# metric="RMSE",
# preProc = c("center", "scale")
```

```
> mod
```

```
Linear Regression | 392 samples | 1 predictor | No pre-  
processing  
Resampling: Leave-One-Out Cross-Validation  
Summary of sample sizes: 391, 391, 391, 391, 391, 391, ...  
Resampling results:
```

RMSE	Rsquared	MAE
429.5254	0.7436498	347.5039

```
> str(model$control$index$Fold001)  
int [1:391] 2 3 4 5 6 7 8 9 10 11 ...
```

```
> plot(mod$pred$obs, type="l");  
  lines(1:392,mod$pred$pred,col="red")
```

