## APRENENTATGE AUTOMÀTIC



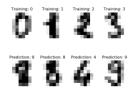
**Objectiu**: aplicar diferents classificadors (regressor logístic, perceptró, arbre de decisió i boscos aleatoris) i entendre les millores d'aplicar kernels, avaluant correctament l'error del model quan és aplicat a dades reals.

**Materials**: els problemes es resoldran mitjançant la llibreria scikit-learn[1] de Python i l'entrega serà un PDF resultat d'un o varis Python Notebooks.

**Puntuació**: Els exercicis s'organitzen en tres nivells de dificultat: A, (sobre 10), dificultat alta; B, (sobre 8), dificultat mitjana i C, (sobre 6), dificultat baixa. Per aprovar els problemes és requisit necessari completar satisfactòriament els problemes de dificultat (C), demostrant així una comprensió fonamental de la matèria.

## **Problema**

S'ha de fer un reconeixedor de números escrits a mà (un OCR de números). Per a fer-ho, s'assumeix feta la part de visió computacional (detecció de text a la imatge, retall de caràcters i normalització de la imatge) i el nostre data set contindrà exemples de números amb un format de imatges de 8 x 8 píxels (que es vectoritza en un vector de 64 posicions amb un recorregut de la imatge per files des de la cantonada superior esquerra) on el valor de cada píxel està dins del rang 0 (negre) – 16 (blanc).



Per tant, fer el reconeixedor de caràcters és el mateix que donada una nova imatge retallada, classificar aquesta imatge en la classe corresponent al número, és a dir, tindrem 10 classes, del "0" al "9".

Per a aquest primer exercici (dificultat C), farem una comparació del rendiment dels següents tres models: el regressor logístic[3], el perceptró[4], l'arbre de decisió[5] i els bosc aleatori [9]. Per a fer l'aprenentatge emprareu el data set "digits"[6] de scikit-learn i la comparació dels models es farà amb el rendiment de cada model amb el conjunt de test que heu triat (expliqueu clarament això) i la matriu de confusió. Compareu i raoneu els resultats trobats.



## APRENENTATGE AUTOMÀTIC

- En el segon exercici (dificultat B), farem la comparació del rendiment dels tres models anteriors amb el data set complert del repositori UCI Machine learning[7]. Primer proveu si canvien els resultats de l'exercici 1 quan s'aplica a noves dades no emprades a aprenentatge ni a test. Després, intentar millorar els resultats re-aprenent els models amb nous conjunts d'aprenentatge, test i validació, emprant el mètode de k-fold crossvalidation[8] per triar el millor model.
- Finalment, en aquest tercer exercici (dificultat A) desenvolupau, al manco, tres característiques de visió per computador i emprau-les per a l'entrament. Analitzau com afecta la seva presència als resultats dels models. Exemple: el nombre de píxels 0 i píxels 1.

Per obtenir les dades amb les que treballareu heu d'aplicar les següents instruccions:

from sklearn.datasets import load\_digits
digits = load\_digits()

**Dates**: El darrer dia per entregar la pràctica serà el dia 16 de gener. La mateixa pràctica es pot entregar posteriorment com a part de l'avaluació complementària, amb una nota màxima de 7, per la dificultat A.

## Referències:

- [1] http://scikit-learn.org/stable/
- [2] https://www.anaconda.com/
- [3] http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.LogisticRegression.html
- [4] http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.Perceptron.html
- [5] https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier
- $\hbox{\it [6]} \ \underline{\text{http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_digits.html}\\$
- [7] https://archive.ics.uci.edu/ml/datasets/optical+recognition+of+handwritten+digits
- [8] http://scikit-learn.org/stable/modules/cross\_validation.html
- [9] https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html