

Welcome to Machine Learning

Konstantin Tretyakov

http://kt.era.ee

IFI Summer School 2014 STACC

Software Technology and Applications Competence Center





Data mining,

Data analysis, Statistical analysis, Pattern discovery, Statistical learning, Machine learning, Predictive analytics, **Business intelligence**, Data-driven statistics, Inductive reasoning, Pattern analysis, Knowledge discovery from databases, Neural networks,















"Unformalizable" problems





"Unformalizable" problems







General rule: IF (X is made of plastic), THEN (X is not edible)

Application in the particular case: X = Rubic's cube

Rubic's cube is not edible

General rule: add(x, y) = function { ... }

Application in the particular case: add(2,4)

General rule: add(x, y) = function { ??? Particular cases: add(2,4) = 6add(5,3) = 8add(1,2) = 3

nduction

Deduction

Given a general rule, make a decision in a particular case

0

Induction

Given particular cases, find a general rule





MNIST dataset

http://yann.lecun.com/exdb/mnist/
Handwritten digits, 28 x 28



MNIST dataset

images = load_images()
labels = load labels()

Let us just use 1000 images training_set = images[0:1000] training labels = labels[0:1000]

MNIST dataset

> training_set[0]

array([0, 0, 0, ..., 254, 241, 198, ...])

> training labels[0]

171



Nearest neighbor method

Ľ

Training set



Nearest neighbor method

def classify(img):
 similarities =
 [similarity(img, p) for p in training_set]
 i = similarities.index(max(similarities))
 return training_labels[i]

Nearest neighbor method

def classify(img):
 similarities =
 [similarity(img, p) for p in training_set]
 i = similarities.index(max(similarities))
 return training_labels[i]

def similarity(img1, img2):
 return -sum(abs(img1 - img2))

Testing the algorithm

test_set = images[1000:2000]
test_labels = labels[1000:2000]

predicted class = [classify(p) for p in test set]

n <mark>succ</mark>esses =

=>

843/1000





clf = KNeighborsClassifier(n_neighbors=1)
clf.fit(training_set, training_labels)

predicted_class = clf.predict(test_set)

=> 855/1000





Some samples may be thrown away







.. or we can SUM instead of MAX


.. or we can SUM instead of MAX



.. or we can SUM instead of MAX



evidence(img, 7) =
 0.0 * similarity(img, tr[0])
 +
 2.0 * similarity(img, tr[3])
 = 20.0





General form

 $K(\boldsymbol{x}_i, \boldsymbol{z})$











 $w_i y_i K(\boldsymbol{x}_i, \boldsymbol{z})$ $f_w(z) =$



How to find the weights?

Find weights, such that the misclassification error rate on the training set is the smallest.

 $f_{w}(z) = \sum w_{i} y_{i} K(x_{i}, z)$

w = argmin_w ErrorRate (f_w, Data)

How to find the weights?

 $f_{w}(z) = \sum w_{i} y_{i} K(x_{i}, z)$

Find weights, such that the approximation to misclassification error on the training set is the smallest.

w = argmin_w Error (f_w, Data)

How to find the weights?

Find weights, such that the error rate on the training set is the smallest + there are many zero weights.

 $f_{w}(z) = \sum w_{i} y_{i} K(x_{i}, z)$

w = argmin_w Error (f_w, Data) + Complexity (w)

Support Vector Machine

from sklearn.svm import SVC

clf = SVC(kernel=`linear')
clf.fit(training set, training labels)

predicted class = clf.predict(test set)

865/1000





Search for the nearest neighbor

def classify(img):

similarities =
 [similarity(img, p) for p in training_set]
 i = similarities.index(max(similarities))
 return training_labels[i]

Search for the nearest neighbor

def classify(img):

similarities =
 [similarity(img, p) for p in training_set]

i = similarities.index(max(similarities))
return training_labels[i]

Inefficient

Search for the nearest neighbor

def classify(img):

nearest_neighbour =
 training_set.find_nearest_neighbour(img)

return nearest neighbour.label

Indexing!











find_nearest_neighbour

if pixel[10,13] > 4:

if pixel[3,24] < 0:

nearest neighbour = A

else:

nearest neighbour = B

else:

nearest neighbour = C

Classification Tree

if pixel[10,13] > 4:

if pixel[3,24] < 0:

class = 1'

else:

class = 2'

else:

class = '3'











General form of a model

 $f_{w}(\boldsymbol{z}) = \sum w_{i} y_{i} K(\boldsymbol{x}_{i}, \boldsymbol{z})$

Search for optimal model parameters





Linear model

f(image) =

pixel I*w / + pixel2*w2 + ... + pixel784*w784

Linear Classification

from sklearn.linear_model import
 LinearRegression,
 LogisticRegression,
 RidgeClassifier,
 LARS,
 ElasticNet,
 SGDClassifier,

> 809/1000




Linear model

f(image) = pixel I*w/ + pixel2*w2 + ... + pixel784*w784







Why should a model, trained on one set of data, work well on future data?

When is it not the case?

How to formalize a new model? How to find parameters?



How to create efficient learning algorithms?





Unsupervised learning Semi-supervised learning On-line learning Active learning Multi-instance learning Reinforcement learning



Tools: R,Weka, RapidMiner, Orange, scikits-learn, MLPy, MDP, PyBrain, ...



"Unformalizable" problems



Deduction and Induction



Model-based

Instance-based

The OCR problem is unusual in that it is ______

The two important perspectives on machine learning are _____-based and _____based.

 $+\lambda$

The "soul" of machine learning is the minimization task

argmin_w

Quiz

Two important components of machine learning:

Quiz

The machine learning algorithms mentioned in this lecture were:

Quiz

Quiz The machine learning algorithms mentioned in this lecture were: K-nearest neighbor classifier **SVM Classification trees** Linear models Neural networks

