# Introduction to Machine Learning

A case based approach

Mannes Poel University of Twente Dept. Computer Science Human Media Interaction



## What is Machine Learning?



**Human Media Interaction** 

**3TU Winterschool 2009** 

## What is Machine Learning?





# What is Machine Learning

- Learning by Machines from Data.
- But there is No Free Lunch, meaning that there is no learning algorithm which performs best on all ML problems.
- A.k.a. Pattern Recognition, Data Mining,



. . . .

# Some applications in Human Media Interaction

**3TU Winterschool 2009** 

- Speech Recognition
- Body:
  - Body detection
  - Pose estimation
  - Gesture recognition
  - Action recognition
- Head:
  - Head detection
  - Head pose estimation

Facial expressions

- Visual Focus of Attention

Main problem for a naturalistic setting is the data:

- How to get naturalistic data?
- Annotations
- Reliability of annotations

**Human Media Interaction** 

# Data in Machine Learning

- It is the interplay between the domain which generated the data and the learning algorithm which determines how good the model is!
- Assumptions on the data:
  - The data is good predictor of the future.
  - The data is or can be made available.
  - The data contains what we want to predict.
- Garbage in Garbage out!!!!!



# Case 1: Detecting heads

#### Challenge

Given a picture or a video stream, can we detect faces in the image (video stream)



# Head detection: Possible approaches

- Shape based:
  - Head as a certain shape, detect such shapes in the picture.
- Skin color or appearance based:
  - Detect skin blobs in the image
  - Largest skin blobs are heads
- Combination of both
- Other

#### Focus: Skin color based approach



# Model: Skin color detector

- Color pictures are coded in RGB values.
- Machine Learning problem is to construct a model *M* which classifies the color of a pixel *x*.
- The possible classes are *skin* or *not skin*.
- Need of data to train (construct) the model *M*.



## Skin color detector

- Use data:
  - Generate skin blob regions from your train data.
  - Determine normalized RGB value for each pixel in the skin blobs, this reduces illumination differences.





#### Visual inspection of skin color data





#### Negative examples



#### Normalized version



**Human Media Interaction** 

#### Data visualization



Observe the different scales of the axis



#### Data visualization: both classes



**Conclusion: skin color classification is non trivial** 



# Histogram based approach

- Given pixel x with value (r,g,b) determine the bin of x. If bin value > t, with t threshold then x is classified as skin.
- Threshold *t* is determined on train set:
  - select *t* with optimal performance on train set.
  - use visualizations
  - this can also be done with the help of Receiver Operator Curve (ROC). Plot Hit rate against False Alarm rate





#### Visualization: data level





## Visualization

#### threshold = 5



raw image



skin pixels



-

skin blobs

raw image

skin pixels







skin blobs

face blob

cropped face







face blob

cropped face





Human Media Interaction

## Visualization

#### threshold = 1



skin pixels

#### skin blobs



face blob









cropped face





raw image



skin pixels



skin blobs

face blob

cropped face









**Human Media Interaction** 

# Curse of dimensionality

In higher dimensions, histogram based approach leads to the so-called "curse of dimensionality"
 x2



number of bins grows exponentially with the dimension of the feature space



3TU Winterschool 2009



 Hence one needs exponentially large quantity of training data. 10 bins per dimension leads to 10<sup>d</sup> bins in dimension d

Solution: Generate models from the data.



#### Generative models

- Generate models M\_skin and M\_non\_skin s.t. M\_skin(x) (M\_non\_skin) is the probability that x is skin (non\_skin).
- x is classified as skin if
  M\_skin(x)>M\_non\_skin(x)
- Examples:
  - Bayesian approach; Bayesian networks
  - Hidden Markov Models



#### Discriminative models

- Generate real valued function f\_skin such that x is classified as skin iff f\_skin(x)>0. Also f\_skin(x) can be seen as a confidence measure.
- Examples:
  - linear discriminant analysis
  - neural networks
  - support vector machines



#### Discrete Discriminative models

- Generate discrete function f\_skin with values in {0,1} such that x is classified as skin iff f\_skin(x)=1
- Examples:
  - Decision Trees





Classification of *x*:  $max_k P(C_k|x) = max_k P(x|C_k)P(C_k)$   $P(C_k)$  can be estimated from trainings set, hence needs model for  $P(x|C_k)$ 



**3TU Winterschool 2009** 

- Need model for *P(x/C);* the probability that
  x is generated by class *C*.
  - for instance a normal distribution.
  - Estimate mean and covariance





**3TU Winterschool 2009** 

 $Cov(x_1, x_2)=0, Var(x_1)=Var(x_2)$ 

 $Cov(x_1, x_2)=0, Var(x_1)>Var(x_2)$ 

x<sub>2</sub>





Equal variances

Variances are different





3TU Winterschool 2009

# Discriminative: Linear classifiers

- A linear classifier partitions in a *clever* way the feature space into two parts using a hyperplane.
- Classification is base on which side of the hyperplane the feature vector lies.
- In 2 dimensions a linear classifier is determined by a line: y=w.x+b. Hence 3 parameters instead of 5 for the Bayesian approach



#### Linear classifier

• Simple decision boundary: hyperplane



# **Decision Trees**

 A decision trees partitions in a *clever* way the feature space into rectangles and assigns a class to ever rectangle.
 Similar to a non-uniform histogram, but the bins are automatically determined by the learning algorithm.





The variables to split upon and the thresholds are determined by the learning algorithm.



## **Neural Networks**

• Structure of simple neurons





 $f(\mathbf{x}) = \phi(\mathbf{w} \bullet \mathbf{x} + \mathbf{b})$ 

33

## **Neural Networks**

• More complex decision boundaries



**Human Media Interaction** 

**3TU Winterschool 2009** 

## Support Vector Machines a.k.a. Kernel Machines

- 1. Map feature vectors to a different space, this determines the kernel.
- 2. Calculate *optimal* linear decision boundary in this transformed space.





# **Overview ML Methods**

Figure from Norbert Jankowski and Krzysztof Grabczewski





#### Case 2: Detecting (upper) bodies



Picasso: Portrait of Wilhelm Uhde

#### Edges play an important role



3TU Winterschool 2009

# Detecting (upper) bodies

- Approach (based on work of Dalal):
  - sliding window
  - window consists of cells
  - for each cell histogram of oriented gradients (edges); edge information
  - feature vector is the concatenation of the cell histograms
  - confidence based classification using Support Vector Machine



#### Detecting (upper) bodies: Example





#### System developed by Ferrari based on the work of Dalal

**3TU Winterschool 2009** 

## Detecting (upper) bodies: Applied to meeting setting





# Case 3: Action Recognition

joined work with Ronald Poppe

- Discriminative approach to action recognition
- Main ingredients
  - Histogram of Silhouette Edges
  - Common Spatial Patterns
  - Simple 1 versus 1 mean CSP based classification technique
- Data set: Blank (Weizmann) data set (<u>http://www.wisdom.weizmann.ac.il/~vision/</u> <u>SpaceTimeActions.html</u>)



# **Action Recognition**





## Weizmann data set

 Video's of 9 persons performing 10 actions: <u>bend</u>, <u>jack</u>, <u>jump</u>, <u>p-jump</u>, <u>run</u>, <u>side</u>, <u>skip</u>, <u>walk</u>, <u>wave1</u> and <u>wave2</u>.

Some actions are performed from left to right by one person and from right to left by another





#### Feature selection process



Region of interest is divided in 4x4 cells. Histogram of silhouette gradients divided into 8 bins, each bin covering 45°. Result: 128-dimensional descriptor per frame.



# Common Spatial Patterns (CSP)

- Frame based feature extraction  $\rightarrow$  sequence of 128-dimensional descriptors.
- How to take into account temporal information?
- CSP is a technique which can be used to discriminate between two classes on bases of the difference in temporal variance.



# Idea behind CSP

Training sequences for two actions *a* and *b*. Each sequence can be seen as  $nxm_p$  with *n* the number of features and  $m_p$  the number of frames.  $C_a(C_b)$  is the concatenation of a the training samples of action a(b).

$$C = C_a C_a^T + C_b C_b^T$$

*C* the temporal variance over time. Assumption: Temporal mean is 0!



# CSP(2)

C is symmetric  $D = UCU^T$ with D a diagonal matrix. Define  $V = \sqrt{D}^{-1}$ then  $I = (VU)C(VU)^T$ 

$$C = C_a C_a^T + C_b C_b^T$$
$$I = (VU)C_a C_a^T (VU)^T + (VU)C_b C_b^T (VU)^T$$

•But this one is also symmetric, hence can diagonalized  $\rightarrow$  W



#### Result:





# CSP(4)

- *L*=*WVU* is a linear transformation such that:
  - the sequences of class *a* have high temporal variance in the first components
  - the sequences of class *b* have low temporal variance in the first components
  - the sequences of class a have low temporal variance in the last components
  - the sequences of class b have high temporal
    variance in the last components



# CSP(5)





3TU Winterschool 2009

# CSP based classification



# CSP based classification



**Human Media Interaction** 

**3TU Winterschool 2009** 

# CSP based classification (2)

•Discriminating function for class *a* and *b*:

$$g_{a,b}(x) = \frac{\|\bar{b} - x'\| - \|\bar{a} - x'\|}{\|\bar{b} - x'\| + \|\bar{a} - x'\|}$$

•Voting function for class a:

$$g_a(x) = \sum_{a \neq b} g_{a,b}(x)$$

Classification of x:  $Cl(x) = argmax_a g_a(x)$ 



# Evaluation (1)

- k=5: this means that every action sequence is reduced to 10-dim. vector
- Leave one person out cross validation
- Performance data set: 95.56%
  - $-\underline{skip}$  of Daria  $\rightarrow$  jumping
  - $-\underline{skip}$  of Ido  $\rightarrow$  running
  - jump of Eli & run of Shahar  $\rightarrow$  walking
  - -<u>wave1</u> of Lyova  $\rightarrow$  wave2



# Evaluation (2)

#### Leave more persons out cross validation

Nr. training persons	Performance
1	64.72%
2	77.82%
3	81.83%
4	84.60%
5	86.63%
6	89.01%
7	91.39%
8	95.56%



# Evaluation

- Divide data set in training, validation and test.
- Models are trained on train set.
- Hyper parameters are determined on validation set.
- Chose best model, including hyper parameters based on performance on validation set.
- Asses performance on test set. This gives the best indication of the performance in "real life"!
- It is not allowed that any information from the test set leaks into the training.



# **Evaluation methods**

- Precision and recall
- Error; percentage of misclassified samples in the test set
- Confusion matrix
- ROC Curve



# **Confusion Matrix**

	Predicted class	
True class	Yes	No
Yes	TP	FN
	True positive	False negative
No	FP	TN
	False positive	True negative

Many valuable statistics, such as error, precision and recall, can be derived from the confusion matrix



# **Concluding remarks**

- For human behavior analysis in a naturalistic setting;
  - annotated data is of utmost importance but hard to get.
  - Uncontrolled environments and context
- In ML there is no "Free Lunch". Best ML approach is determined by the problem under consideration.



# More information

- Video lectures:
  - Introduction to Machine Learning
  - Machine Learning in Vision
- Books:
  - E. Alpaydin: Introduction to Machine Learning
  - C.M. Bishop: Pattern Recognition and Machine Learning
- Journals:
  - Machine Learning

Im - Journal of Machine Learning

3TU Winterschool 2009





## Case 4: Visual Focuses of Attention

Challenge

# Can we build a model which determines the head of a person and where the person is looking at.



# VFoA recognition system

work of Elisa Ricci @ Idiap

- VFoA (*who is looking at whom or what*) is defined by eye gaze
- Eye gaze is difficult to infer: head pose is used as an indication of VFoA



- IDIAP system for VFoA recognition
  - Fully automatic
  - Running close to realtime
  - Handling faces at low/mid resolution





# Head tracking and pose estimation

- Goal: provide as outputs the location of the head and the head pose (pan/tilt/roll angles).
- *Joint* head tracking and pose estimation in a Bayesian framework
  - Mixed state particle filter: the state X contains both continuous and discrete variables.
  - Automatic *initialization* with a multiview face detector.
  - Exemplar based tracking: a reference models for each possible pose.
  - Reference models *learned* offline with a *discriminative approach* from the PRIMA-POINTING database.



 $X = (S, \theta, r)$ 

S: head location and size  $\theta$ : out-plane rotations r: in-plane rotations





# Head tracking and pose estimation

- Observation model
  - $Y = (Y^{text}, Y^{col})$ : independent features for textures and color
  - Features are computed with integral images
  - Y<sup>text</sup>: concatention of Edge Orientation Histograms (EOH)
  - Y<sup>col</sup>: binary masks extracted by a gaussian skin color model learned offline and adapted online
- Likelihood function



$$= e^{-\lambda_{text}D_{text}(Y^{text}(S_i),R^{text}(\theta,r))} e^{-\lambda_{col}D_{col}(Y^{col}(S_i),R^{col}(\theta,r))}$$
  
Texture likelihood Color likelihood



3TU Winterschool 2009

 $p(Y \mid X = (S, \theta, r)) = p_{text}(Y^{text}(S) \mid \theta, r) p_{col}(Y^{col} \mid \theta, r)$ 

# **VFoA** recognition

- Gaussian distribution models to estimate for a given person the VFOA  $f^t$  given the head pose  $\theta^t$  (using the mean of the particle distribution)
  - Gaussian distribution for regular target

$$p(\theta^{t} \mid f^{t} = T_{k}) = N(\theta^{t}; \mu_{k}, \Sigma_{k})$$

Uniform distribution for unfocused



# VFoA recognition

- Representation of the head pose distribution only by its mean value implies a loss of information
- Exponential distribution models

$$p(\theta^{t} \mid f^{t} = T_{k}) = \lambda e^{-\lambda \rho \left(\pi_{k}^{t}(\theta), \Pi_{k}^{t}\right)}$$

 $\rho(\pi_k^t(\theta), \Pi_k^t)$  distance between distributions

$$\Pi_{k}^{t} \left( \theta = l \right) = \frac{N \left( l \; ; \mu_{k}, \Sigma_{k}^{t} \right)}{\sum_{l'} N \left( l'; \mu_{k}, \Sigma_{k}^{t} \right)}$$

pdf modeling focusing at target k





3TU Winterschool 2009

# **VFoA** recognition





# My Research Interests

- Applied Machine Learning:
  - Vision based Human Behavior Computing:
    - Multi Modal Laughter Recognition (MLMI 2008)
    - Pose Estimation (AMDO 2008)
    - Action Recognition (FG 2008)
  - Brain Computer Interfacing
    - BCI as one of the Interaction Modalities for Games
    - Challenge: dealing with artifacts



# Content

- Introduction to Machine Learning (ML)
- Some cases:
  - Head detection
  - Upper Body Detection
  - Action Recognition
  - Visual Focus of Attention
- Evaluation
- Some concluding remarks

