

# 9.913 Pattern Recognition for Vision

Class I - Overview

*Instructors: B. Heisele, Y. Ivanov, T. Poggio*

# TOC

- Administrivia
- Problems of Computer Vision and Pattern Recognition
- Overview of classes
- Quick review of Matlab

# Administrivia

- **Instructors:**
  - Bernd Heisele
  - Yuri Ivanov
  - Tomaso Poggio
  
- **Meet**
  - Th 10-12
  
- **Credits: 10H**
  
- **Assignments:**
  - Small weekly
  - Paper presentation/discussion
  - Final project

# Syllabus

<i>Sep 9 - Overview, Introduction</i>	<i>Y&amp;B</i>
<i>Sep 16 - Vision - Image formation and processing</i>	<i>Y</i>
<i>Sep 23 - Vision – Feature extraction I</i>	<i>B</i>
<i>Sep 30 - PR/Vis - Feature Extraction II/Bayesian decisions</i>	<i>B&amp;Y</i>
<i>Oct 7 - PR - Density estimation</i>	<i>Y papers</i>
<i>Oct 14 - PR – Clasification</i>	<i>B</i>
<i>Oct 21 - Biological Object Recognition</i>	<i>T</i>
<i>Oct 28 - PR - Clustering</i>	<i>Y&amp;B proj</i>
<i>Nov 4 - Paper Discussion</i>	<i>All</i>
<i>Nov 11 - App I - Object Detection/Recognition</i>	<i>B</i>
<i>Nov 18 - App II - Morphable models</i>	<i>T&amp;B</i>
<i>Nov 25 - No class - Thanksgiving day</i>	
<i>Dec 2 - App III - Tracking</i>	<i>C&amp;Y</i>
<i>Dec 9 - App IV - Gesture and Action Recognition</i>	<i>Y</i>
<i>Dec 16 - Project presentation</i>	<i>All</i>

# Course Materials

- Books
  - Duda, Hart and Stork, *Pattern Recognition*
  - Optional - Mallot, *Computational Vision: Information Processing in Perception and Visual Behavior*
  - *Suggested further reading*
    - *Vision: Forsyth, Ponce, “Computer Vision: a Modern Approach”*
    - *Machine Learning: Hastie, Tibshirani, Friedman, “The Elements of Statistical Learning”*
- Slides, links
- Papers, notes, tutorials
- Office hours
  - No set hours – e-mail, or call

# Computer Vision

## Problems of Computer Vision

- Shape from Shading
- Stereo
- Structure from Motion
- Tracking
- Object Detection/Recognition
- Activity Detection/Recognition
- ... many more....

# Machine Vision - Applications

- Automatic quality control
- Robotics
- Perceptual interfaces – human/machine interactions
- Surveillance
- ...

Images removed due to copyright considerations.

# Computer Vision – Shape-from-Shading

Find 3D surface parameters from a single image based on shading

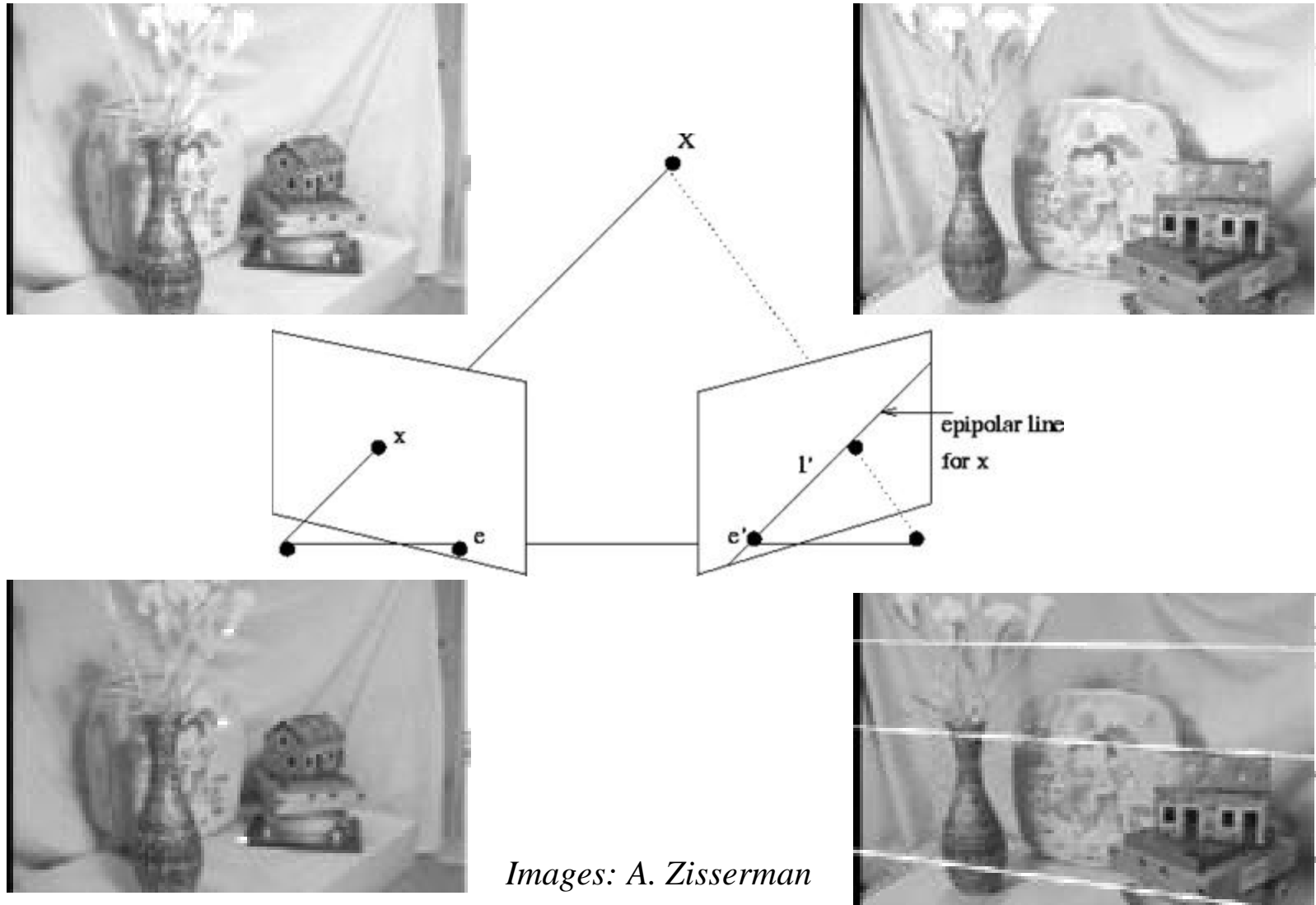
Images removed due to copyright considerations.

*Images: NASA*



# Computer Vision - Stereo

Reconstruct 3D surface from 2 2D images taken simultaneously



*Images: A. Zisserman*

# Computer Vision - Structure-from-Motion

Reconstruct 3D surface from 2D images taken from a moving camera

*Images* \_\_\_\_\_

Photos removed due to copyright considerations.  
Please see Figure 9 in: Azarbayejani, and Pentland.  
"Recursive Estimation of Motion, Structure, and Focal  
Length." IEEE Transactions on Pattern Analysis and  
Machine Intelligence 7, no. 6 (June 1995): 562-575.

*Structure* \_\_\_\_\_

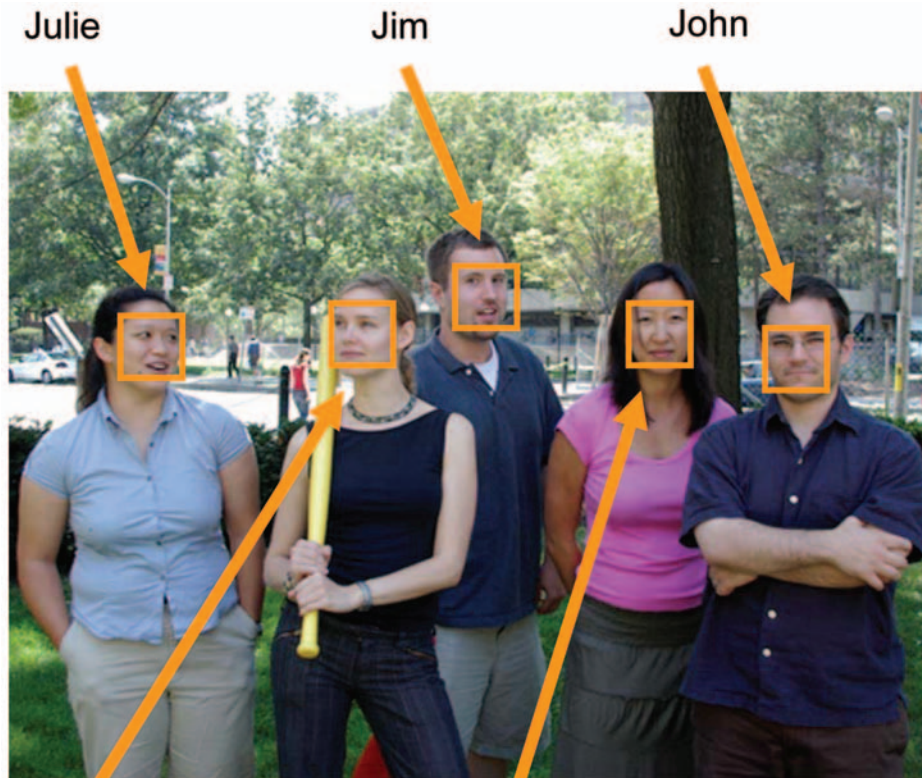
*Texture* \_\_\_\_\_

*Textured model* \_\_\_\_\_

# Computer Vision - Object Detection \ Recognition

Find objects in the image, determine what they are

Eg: Face detection and recognition:



Julie

Jim

John

Jane

Jessica

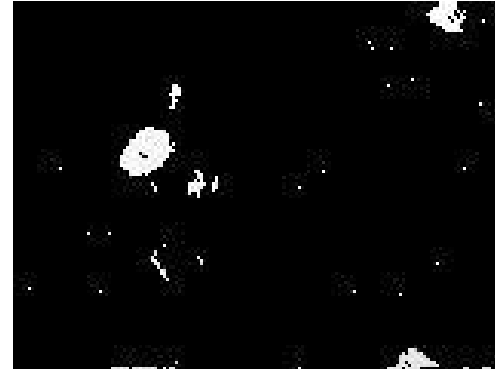
Photograph by MIT OCW.

# Computer Vision - Tracking

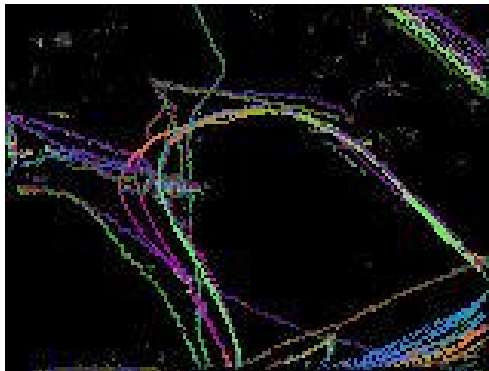
Determine objects' positions over multiple frames



Camera view



Connected components



Trajectories over time



An object

*Tracker – Chris Stauffer, Eric Grimson*

Figures and photographs from: Stauffer, and Grimson. "Learning patterns of activity using real-time tracking." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22, no. 8 (August 2000): 747-757. Courtesy of IEEE, Chris Stauffer, and Eric Grimson. Copyright 2000 IEEE. Used with Permission.

# Pattern Recognition

## “Big Four” Problems of Machine Learning

- Classification
- Density Estimation
- Clustering
- Regression

# Classification

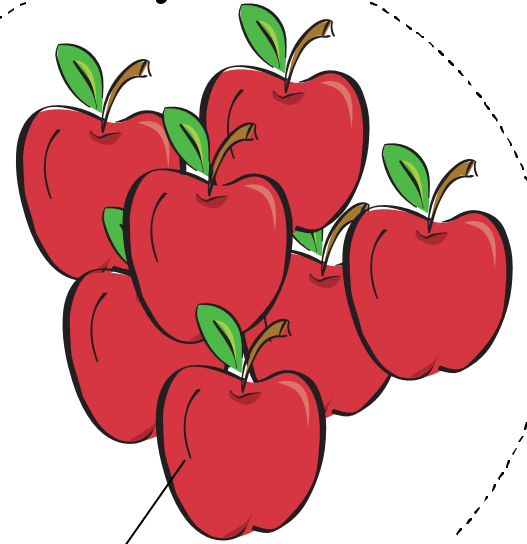
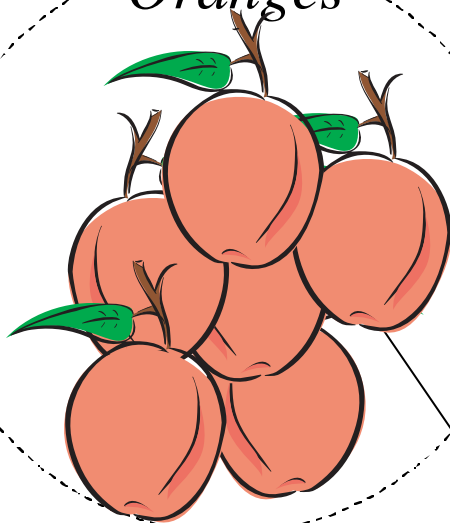
Test data

Decision boundary

Class model



*Oranges*



*Apples*

Training data

Images by MIT OCW.

# Density Estimation

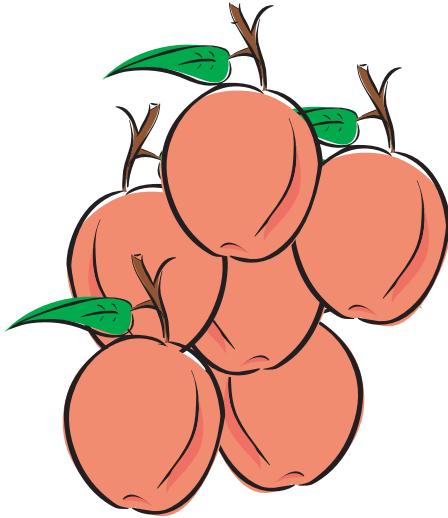
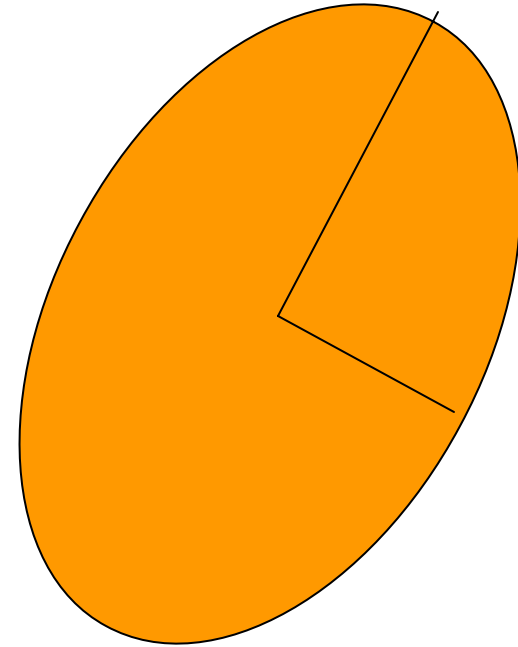
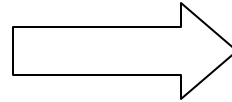


Image by MIT OCW.

Individual samples

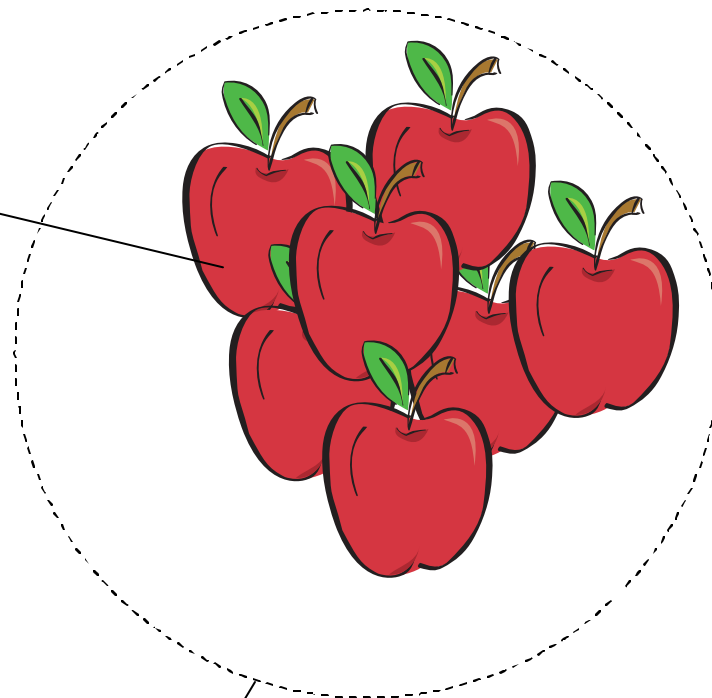
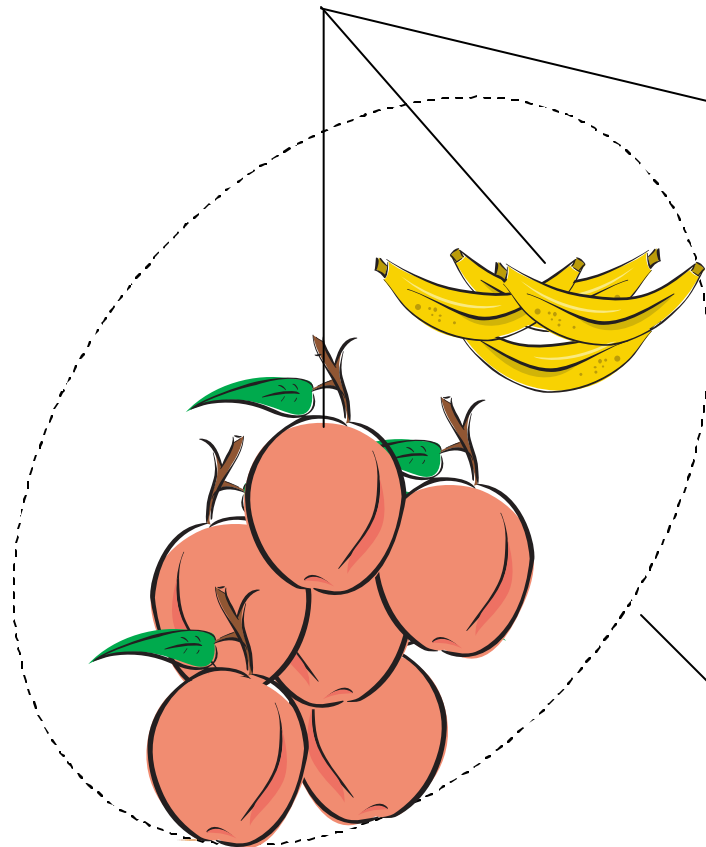


$\mu, \Sigma$

Generating density

# Clustering

Data

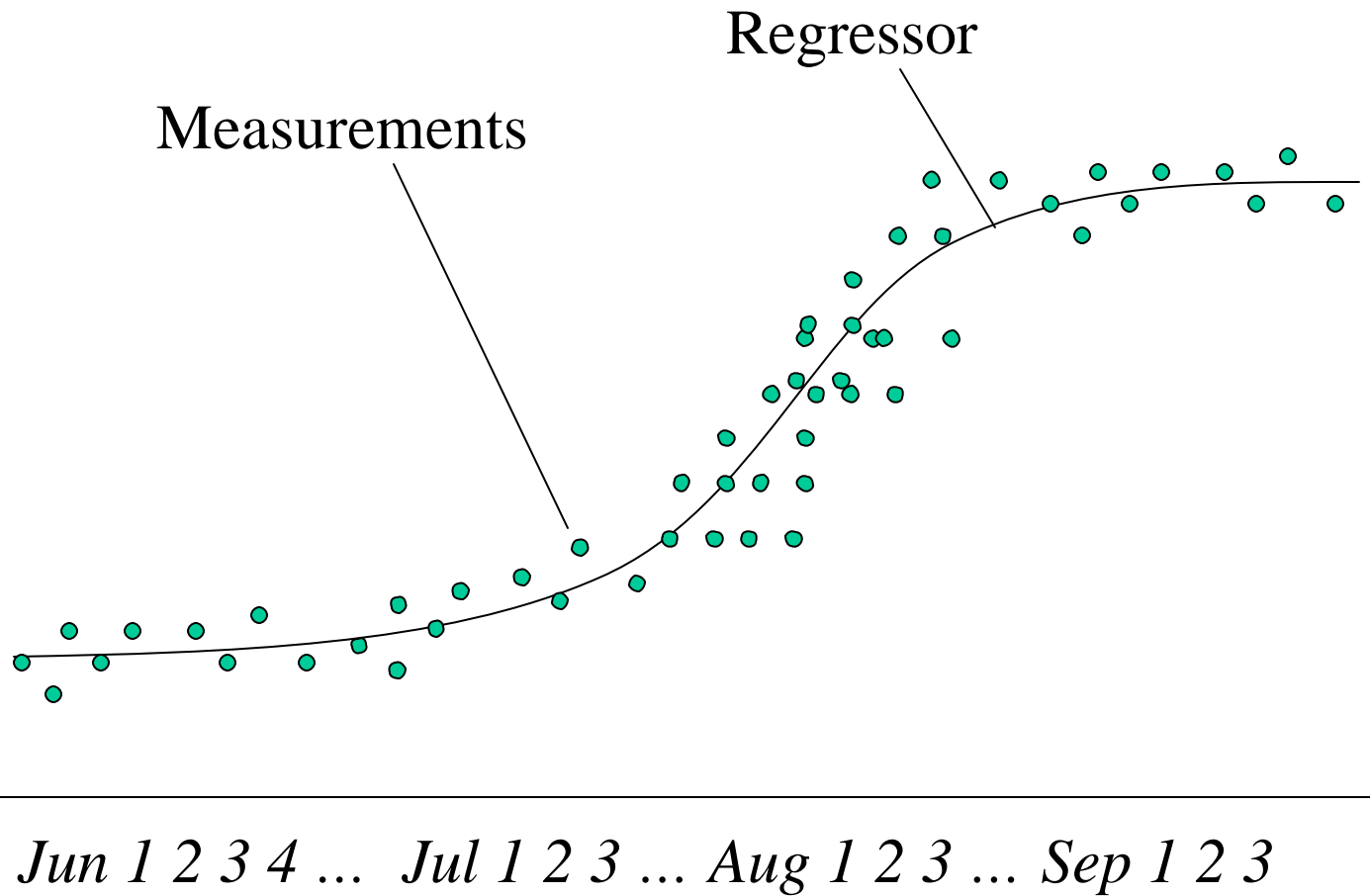
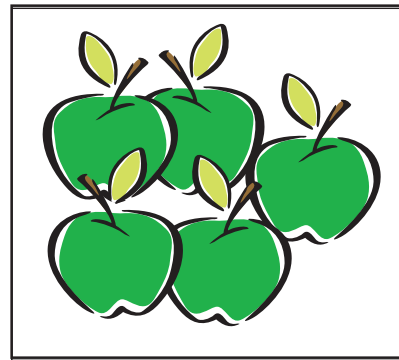
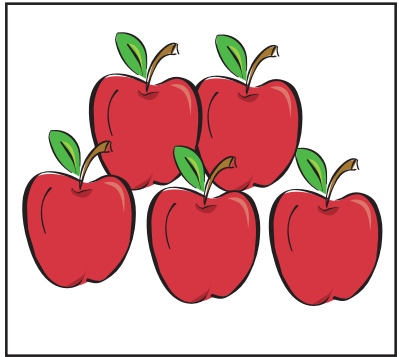


Spatial clusters

Images by MIT OCW.



# Regression



Images by MIT OCW.

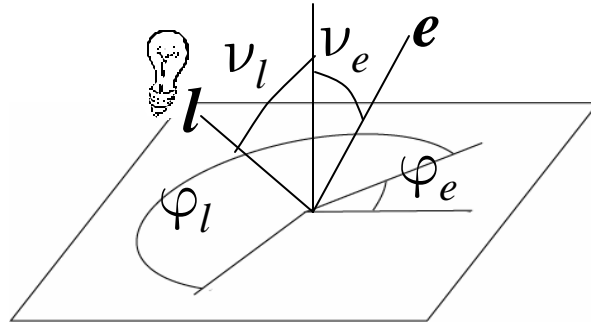
# Vision – Image Formation and Processing

## Topics

- Image formation
- Color spaces
- Point operators
- Neighborhood operators
- Edge detection
- Motion

# Vision - Fundamentals

- Image formation



- Source properties:

- Type of source
- Occlusions / shadows
- Direction

- Surface properties:

- Albedo (fraction of light reflected),
- Shape,
- Smoothness,
- Orientation

- Imager:

- Color properties
- Distortion
- Focal length
- F-stop (how wide the iris)
- Orientation

- Representation:

- Bit depth
- Color space

# Vision - Color

- Light is perceived by two types of receptors in the human eye
  - Rods
  - Cones
- Color is perceived by 3 types of cones
  - “Red” (64%)
  - “Green” (34%)
  - “Blue” (2%)

120,000,000 rods  
6,000,000 cones

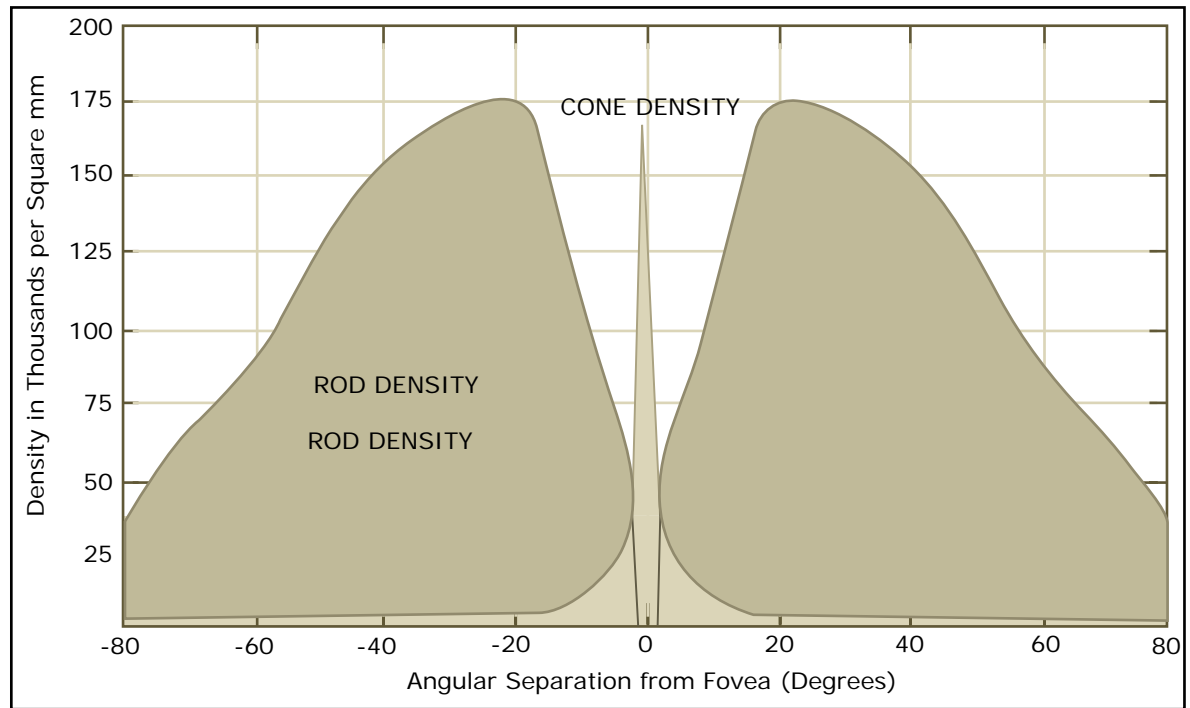


Figure by MIT OCW.

# Color – CIE Chromaticity Diagram

All Visible Colors (CIE 1931 Chromaticity diagram):

Figure removed due to copyright reasons. Please see: <http://www.cie.co.at/cie/>  
and  
<http://www.ledtronics.com/datasheets/Pages/chromaticity/097b.htm>

# Point Operators – Image Arithmetic

Point operators query and affect only a single pixel at a time

Thresholding of a very famous person:

Original photo and accompanying processed photo both removed due to copyright reasons.

## Topics

- Intro to Pattern Recognition and Machine Learning
- Bayes rule
- Normal density
- Minimum error rate classification
- Decision surfaces
- ROC curves and classifier performance

# Bayesian Decision Theory

Bayes Rule:

$$\begin{array}{c} \textit{Posterior} \\ \diagdown \\ \boxed{P(\mathbf{w} | x)} \end{array} = \frac{\begin{array}{c} \textit{Likelihood} \\ \diagdown \\ \boxed{P(x | \mathbf{w})} \end{array} \begin{array}{c} \textit{Prior} \\ \diagdown \\ \boxed{P(\mathbf{w})} \end{array}}{\begin{array}{c} \boxed{P(x)} \\ \diagup \\ \textit{Evidence} \end{array}}$$

Image removed due to copyright considerations.

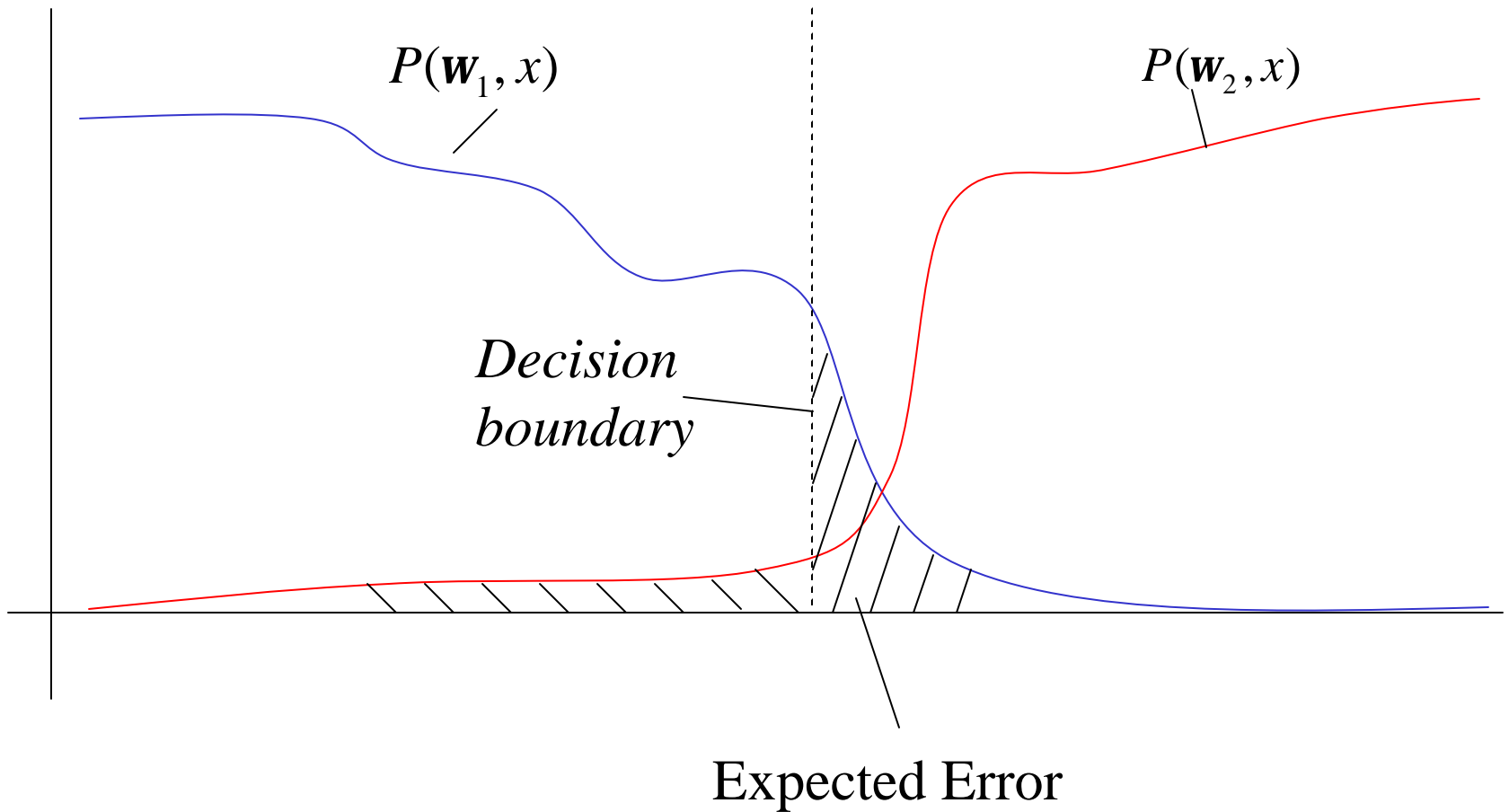
*Rev. Thomas Bayes  
(1702-1761)*

*Fundamental Law of Probability and Statistics*



# Expected Error

How do we compare two class models given the data?



# Likelihood Ratio Test

How do we compare two class models given the data?

$$\frac{P(\mathbf{w}_1 | x)}{P(\mathbf{w}_2 | x)} \lessgtr 1 \Rightarrow$$

$$\frac{P(x | \mathbf{w}_1)P(\mathbf{w}_1)}{P(x | \mathbf{w}_2)P(\mathbf{w}_2)} \lessgtr 1 \Rightarrow$$

$$\frac{P(x | \mathbf{w}_1)}{P(x | \mathbf{w}_2)} \lessgtr \frac{P(\mathbf{w}_2)}{P(\mathbf{w}_1)}$$

|  
Likelihood Ratio Test

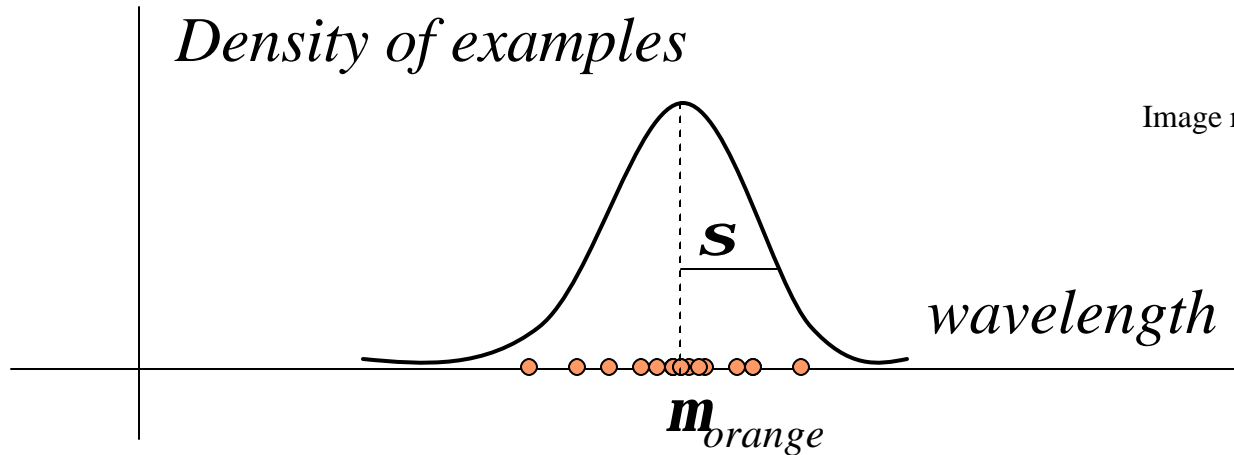
## Topics

- Parametric density estimation
  - ML parameter estimation
  - Bayesian parameter estimation
- Non-parametric
  - K-Nearest Neighbors
  - Parzen Windows

# Density Estimation

Say we need to teach a robot about oranges

Main feature: color



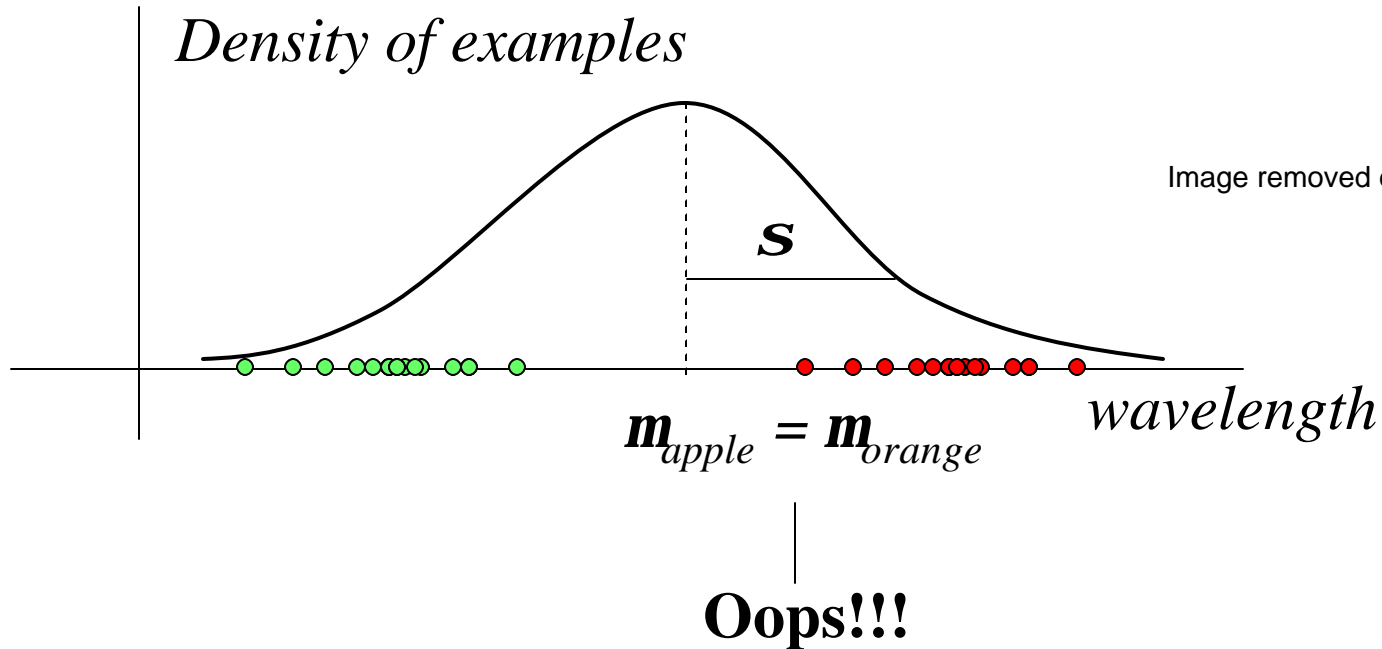
Assume a Gaussian:  $P(w_i | \mathbf{m}, \mathbf{s}) = \frac{1}{\sqrt{2\pi s}} e^{-\frac{1}{2} \left( \frac{w_i - \mathbf{m}}{s} \right)^2}$

Need  $\mu$  and  $\sigma$  to maximize  $P(w_1 w_2 w_3 \dots w_N / \mathbf{m}, \mathbf{s})$

*Solution: factor, log, differentiate and set to 0*

# Density Estimation

Say we need to teach a robot about apples



Need a multi-modal distribution, or a *non-parametric* method

# Pattern Recognition – Clustering

## Topics

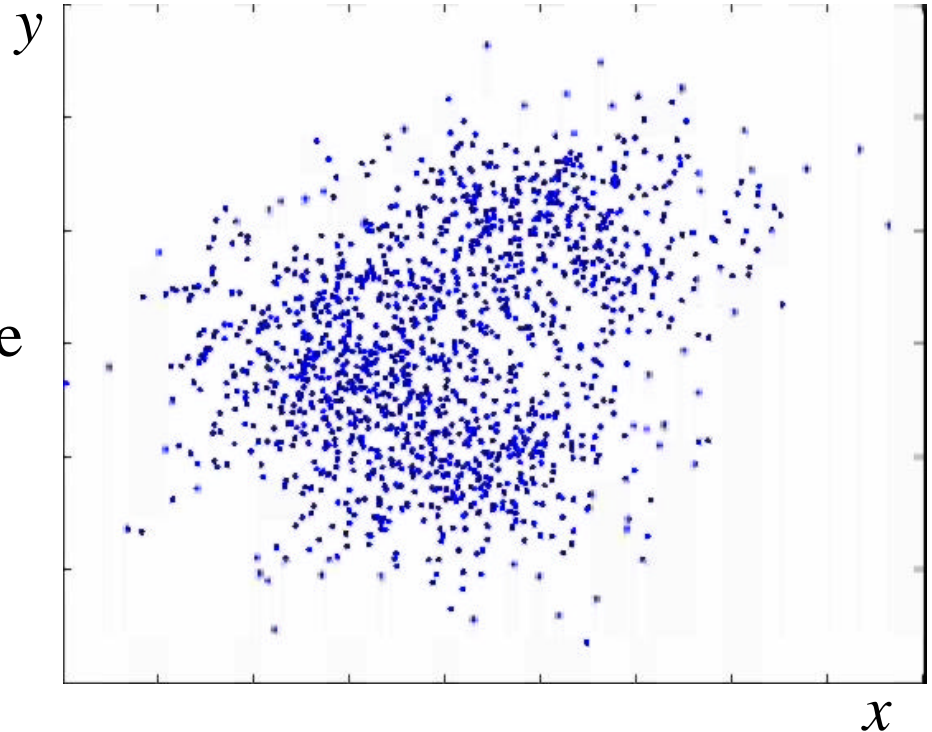
- K-Means algorithm
- EM algorithm
- Trees
- Clustering for tracking

# Clustering/Vector Quantization

Find tight groups in the data

*K-Means* algorithm

1. Randomly place  $K$  centers
2. Assign points to the closest one
3. Compute new centers (means)
4. Go to 2 until convergence



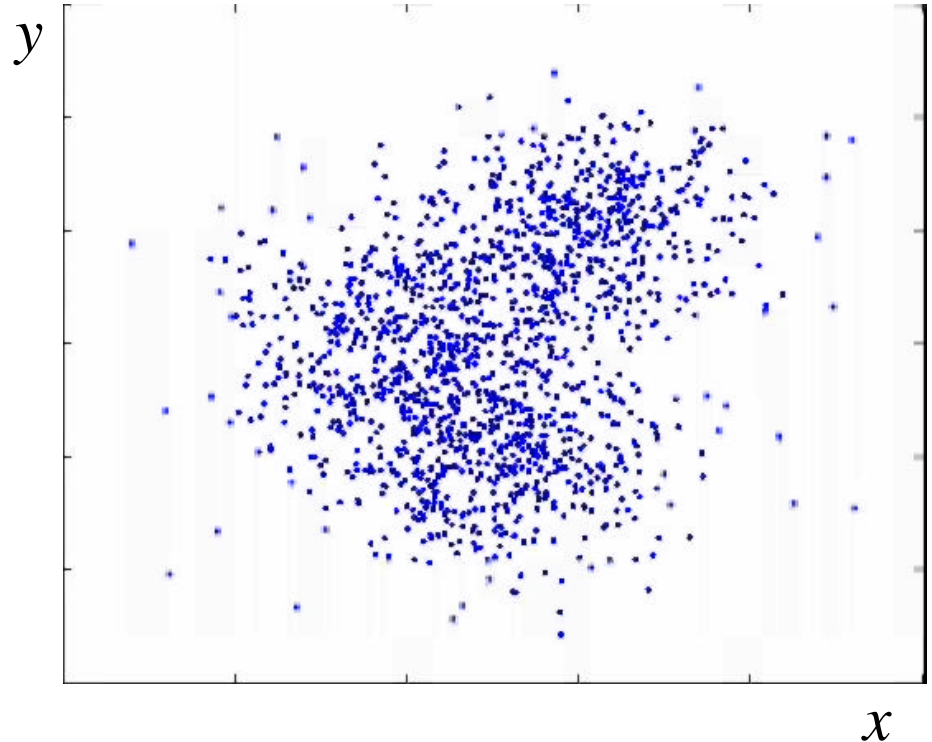
“Hard” label assignment

# Clustering/Vector Quantization

Find tight groups in the data

*EM* algorithm

1. Randomly place  $K$  Gaussians
2. Compute the *posterior* for each point
3. Average data with respect to posterior of each class
4. Go to 2 until convergence

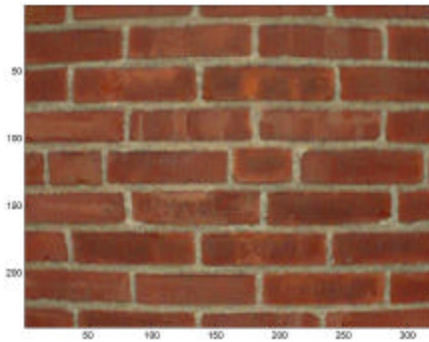


“Soft” label assignment

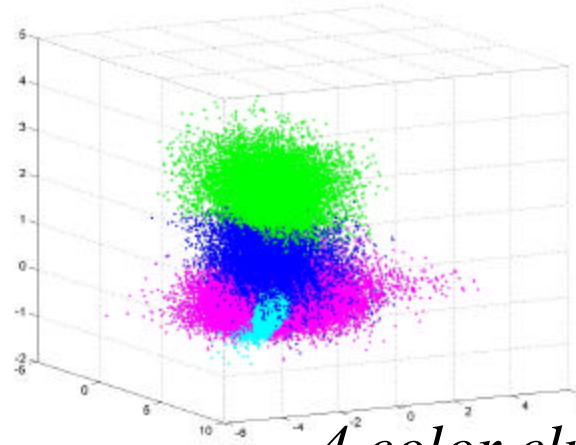
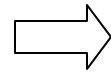


# Clustering – Another Example

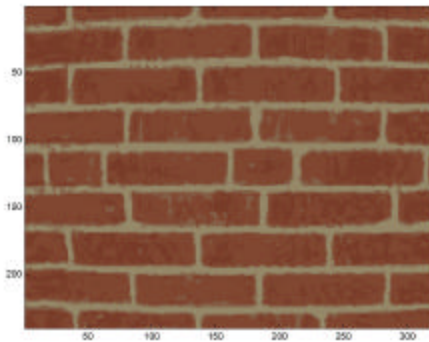
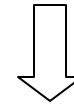
In the image we can cluster pixels in the RGB space



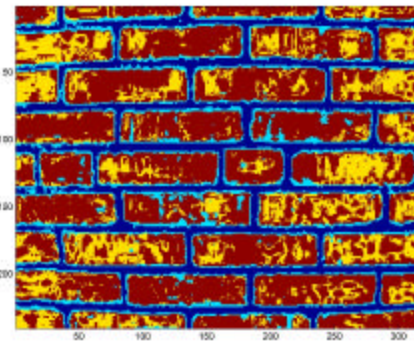
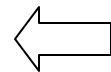
*Source – 24 bit/pix*



*4 color clusters*



*Color-quantized – 2 bit/pix*



*Cluster assignments*

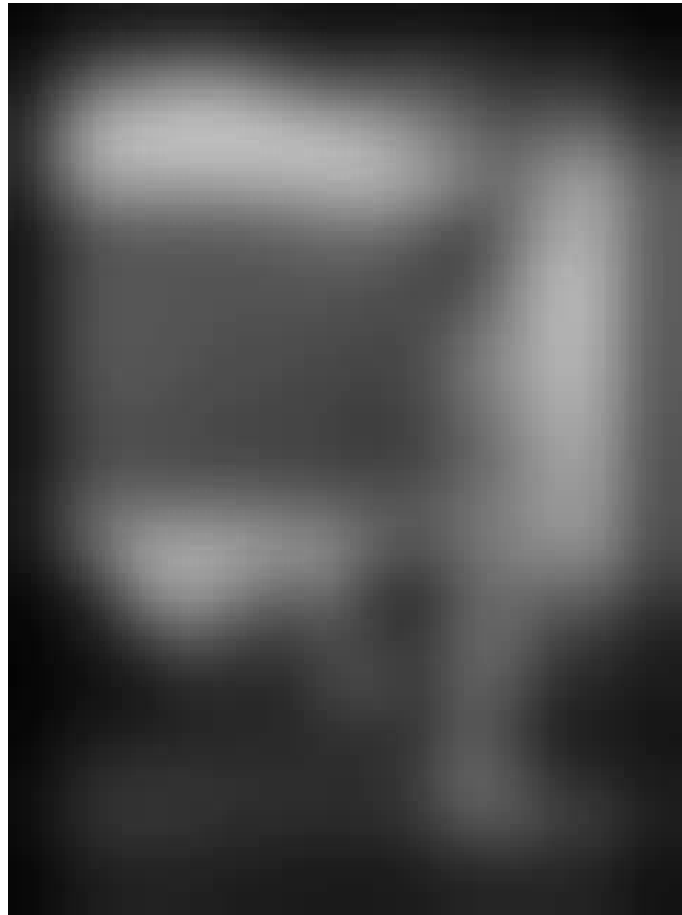
# Vision – Movement, Activity, Action

## Topics

- Motion
- Gesture
- Surveillance

# Application: Motion and Gesture Recognition

Is motion important for object recognition?



Photos and figures from: Bobick, A., and J. Davis. "The Representation and Recognition of Action Using Temporal Templates." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23, no. 3 (2002). Courtesy of IEEE, A. Bobick, and J. Davis. Copyright 2002 IEEE. Used with Permission.

# Motion

One way to do it – turn it into object recognition:

Video:



Cumulative  
MEI:



Now shape can be matched

Problems: Direction, Segmentation, Only appropriate for determining gross body motion...

Photos and figures from: Bobick, A., and J. Davis. "The Representation and Recognition of Action Using Temporal Templates." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23, no. 3 (2002). Courtesy of IEEE, A. Bobick, and J. Davis. Copyright 2002 IEEE. Used with Permission.

## Linear Dynamic Systems:

- Systems where the next **state** is a linear combination of the the previous state, a control signal, and noise.
- **Observations** are a linear combination of the current state and noise

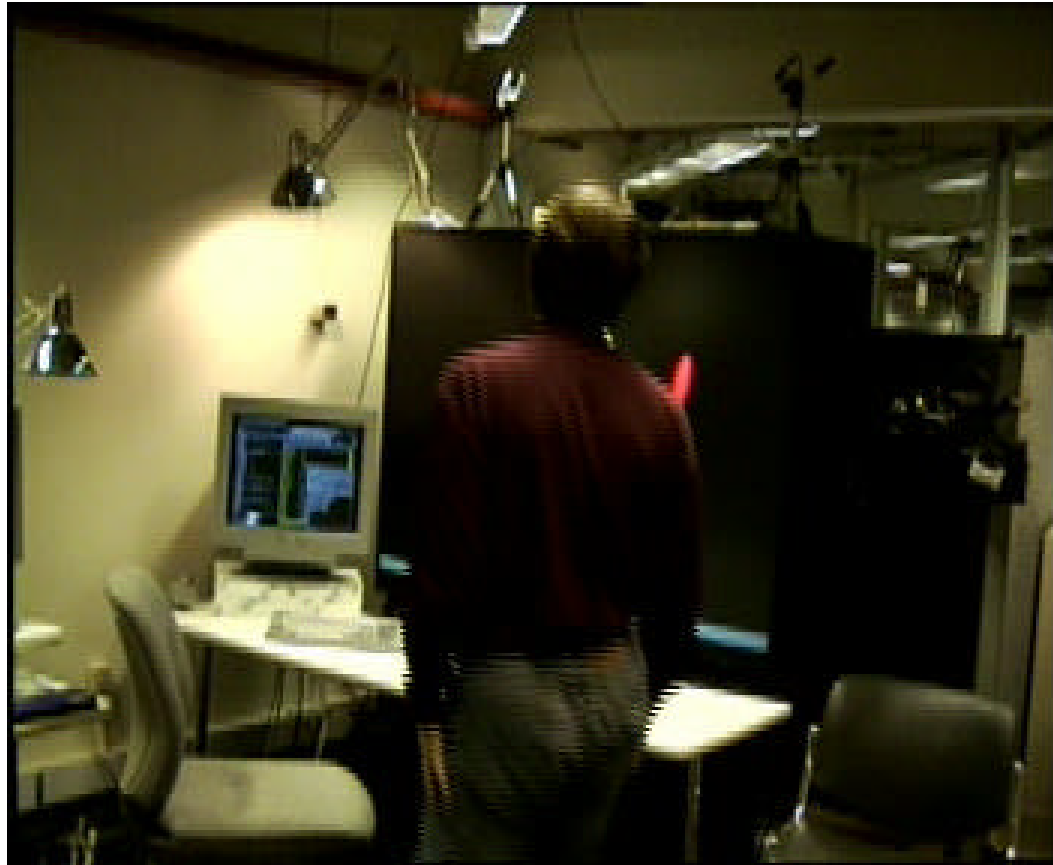
$$\mathbf{x}_{t+1} = \mathbf{\Phi}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{L}_t \boldsymbol{\xi}_t$$

$$\mathbf{y}_{t+1} = \mathbf{H}_t \mathbf{x}_{t+1} + \boldsymbol{\theta}_t$$

# Tracking

- Trackers combine system models  $\{F, B, H\}$ , with observations to estimate the state.
- *Unimodal* trackers, like Kalman Filters, maintain a single estimate of the state
- *Miltimodal* tracker frameworks combine multiple unimodal trackers to track multimodal distributions. (e.g.: Multiple Hypothesis Testing, Particle Filtering)

# Tracking - Example



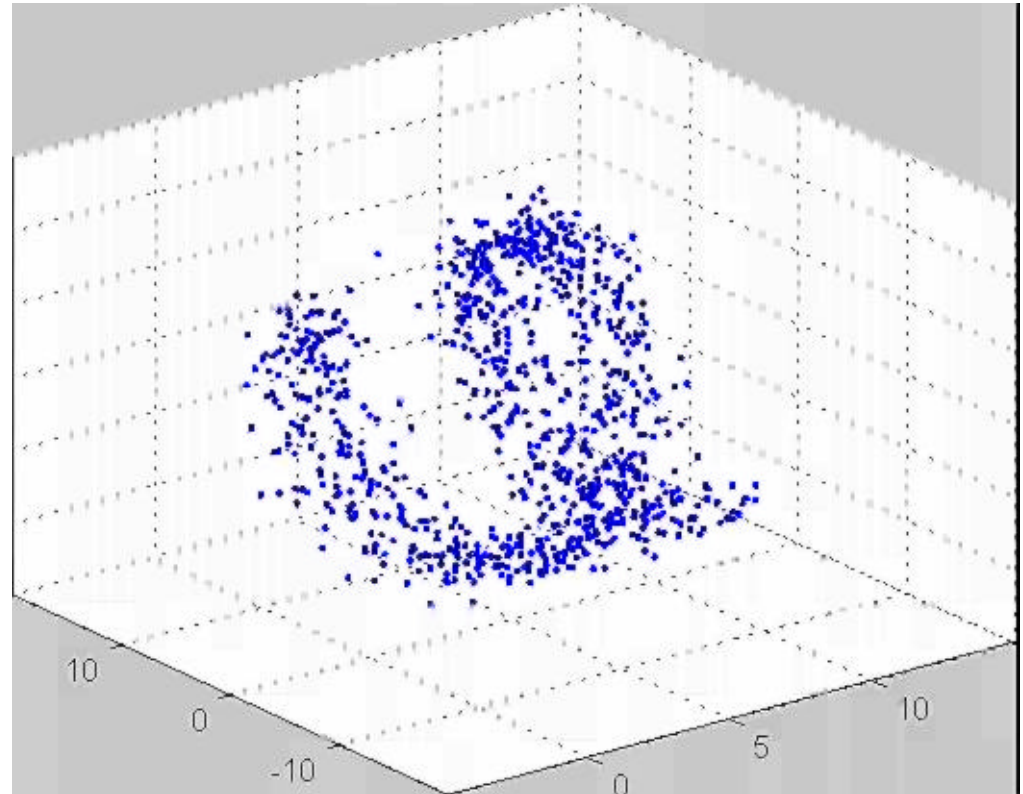
Courtesy of Chris Wren. Used with permission.

# Gesture

Tracking allows modeling gesture as time series:

*Hidden Markov Model (HMM):*

- A mixture of simple densities
- With dynamic constraints



After this sequential density is estimated things are easy again



# Applications IV - Surveillance



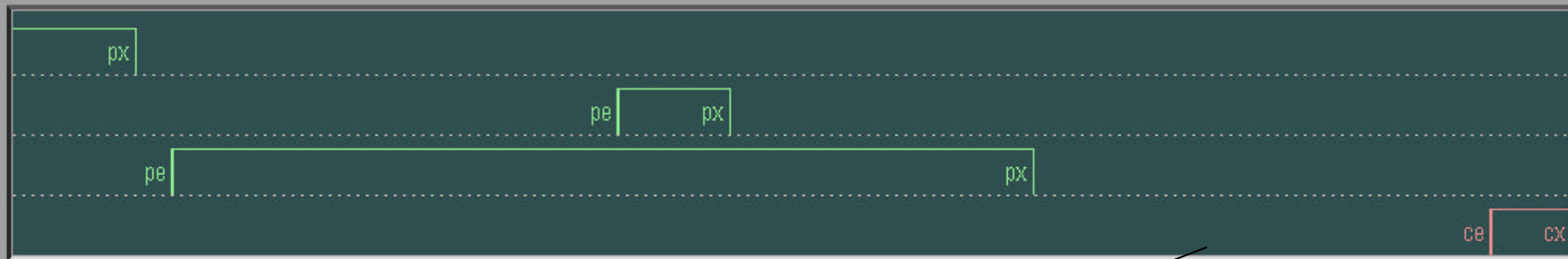
# Parse: Car-Pass-Through

Action label

Component labels

Object tracks

```
Track-12: person_enter(543) SKIP person_exit(543) P = 0.05081913  
Event 33: 0.175000 2 33  
Segmentation:  
Person pass-through, frames 1871 - 1889  
P = 0.17499998  
Track-13: person_enter(707) person_exit(707) P = 0.17499998  
Event 36: 0.174871 2 36  
Segmentation:  
Person pass-through, frames 1799 - 1938  
P = 0.17481139  
Track-14: person_enter(665) person_exit(665) P = 0.17481139  
Event 41: 0.040299 3 41  
Segmentation:  
Car pass-through, frames 2012 - 2025  
P = 0.04029916  
Track-15: car_enter(790) SKIP car_exit(790) P = 0.04029918
```



Temporal extent

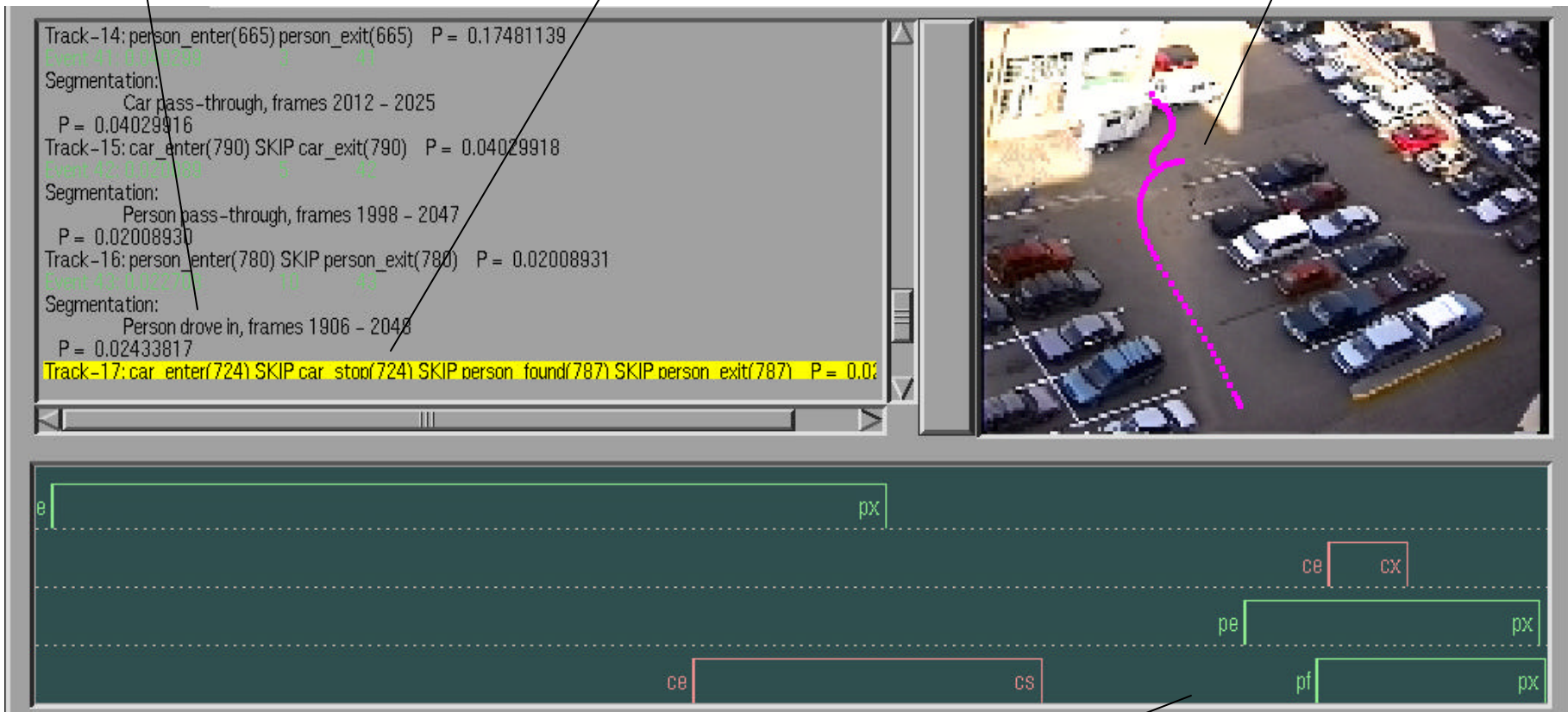
Figures from: Ivanov, Yuri, Chris Stauffer, Aaron Bobick, and W. E. L. Grimson. "Video Surveillance of Interactions." *IEEE Workshop on Visual Surveillance (ICCV 2001)* (1999). Courtesy of IEEE, Y. Ivanov, C. Stauffer, A. Bobick, W. E. L. Grimson. Used with Permission.

# Parse: Drive-In

Action label

Component labels

Object tracks



Temporal extent

Figures from: Ivanov, Yuri, Chris Stauffer, Aaron Bobick, and W. E. L. Grimson. "Video Surveillance of Interactions." *IEEE Workshop on Visual Surveillance (ICCV 2001)* (1999). Courtesy of IEEE, Y. Ivanov, C. Stauffer, A. Bobick, W. E. L. Grimson. Copyright 1999 IEEE. Used with Permission.

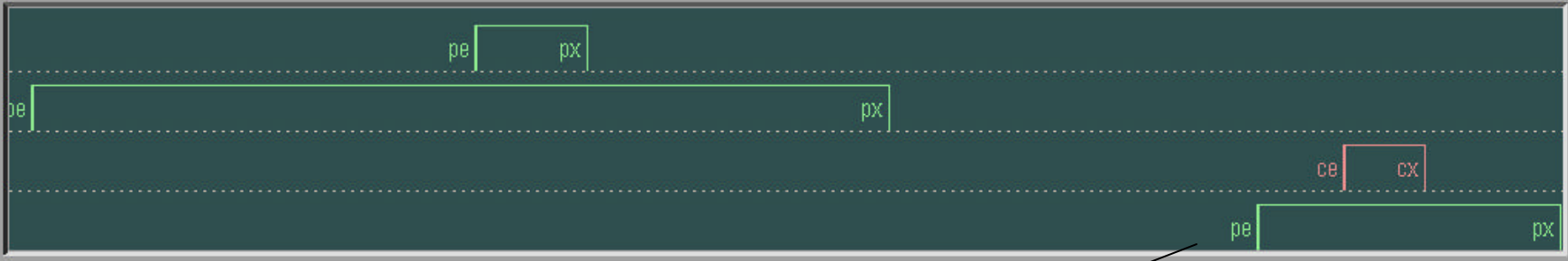
# Parse: Person-Pass-Through

Action label

Component labels

Object tracks

```
Track-13: person_enter(707) person_exit(707) P = 0.17499998  
Event 36: 0.174811 2 36  
Segmentation:  
Person pass-through, frames 1799 - 1938  
P = 0.17481139  
Track-14: person_enter(665) person_exit(665) P = 0.17481139  
Event 41: 0.040299 3 41  
Segmentation:  
Car pass-through, frames 2012 - 2025  
P = 0.04029916  
Track-15: car_enter(790) SKIP car_exit(790) P = 0.04029918  
Event 42: 0.020089 5 42  
Segmentation:  
Person pass-through, frames 1998 - 2047  
P = 0.02008930  
Track-16: person_enter(780) SKIP person_exit(780) P = 0.02008931
```



Temporal extent

Figures from: Ivanov, Yuri, Chris Stauffer, Aaron Bobick, and W. E. L. Grimson. "Video Surveillance of Interactions." *IEEE Workshop on Visual Surveillance (ICCV 2001)* (1999). Courtesy of IEEE, Y. Ivanov, C. Stauffer, A. Bobick, W. E. L. Grimson. Copyright 1999 IEEE. Used with Permission.

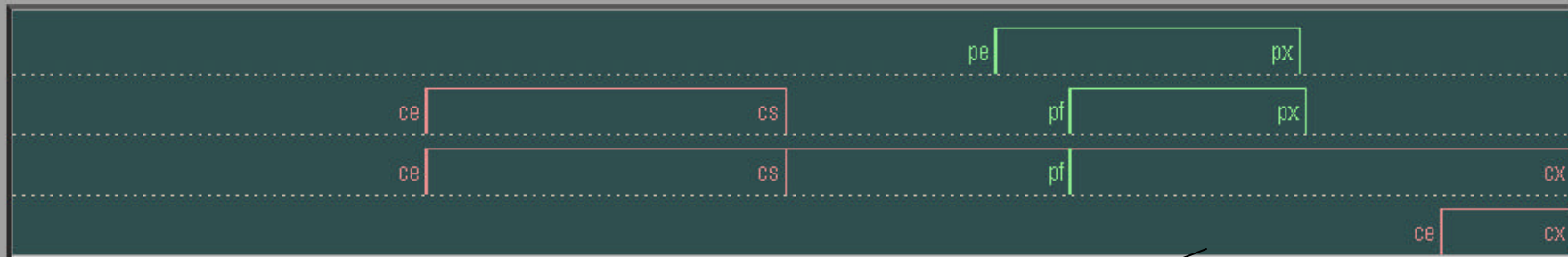
# Parse: Car-Pass-Through

Action label

Component labels

Object tracks

```
Track-16: person_enter(780) SKIP person_exit(780) P = 0.02008931  
Event 43: 0.017758 10 43  
Segmentation:  
Person drove in, frames 1906 - 2048  
P = 0.02433817  
Track-17: car_enter(724) SKIP car_stop(724) SKIP person_found(787) SKIP person_exit(787) P = 0.00117883  
Event 46: 0.017783 13 46  
Segmentation:  
Person drop off, frames 1906 - 2091  
P = 0.01780757  
Track-18: car_enter(724) SKIP car_stop(724) SKIP person_found(787) SKIP car_exit(724) P = 0.0168  
Event 47: 0.050539 3 47  
Segmentation:  
Car pass-through, frames 2070 - 2091  
P = 0.05053889  
Track-19: car_enter(816) SKIP car_exit(816) P = 0.05053889
```



Temporal extent

Figures from: Ivanov, Yuri, Chris Stauffer, Aaron Bobick, and W. E. L. Grimson. "Video Surveillance of Interactions." *IEEE Workshop on Visual Surveillance (ICCV 2001)* (1999). Courtesy of IEEE, Y. Ivanov, C. Stauffer, A. Bobick, W. E. L. Grimson. Copyright 1999 IEEE. Used with Permission.

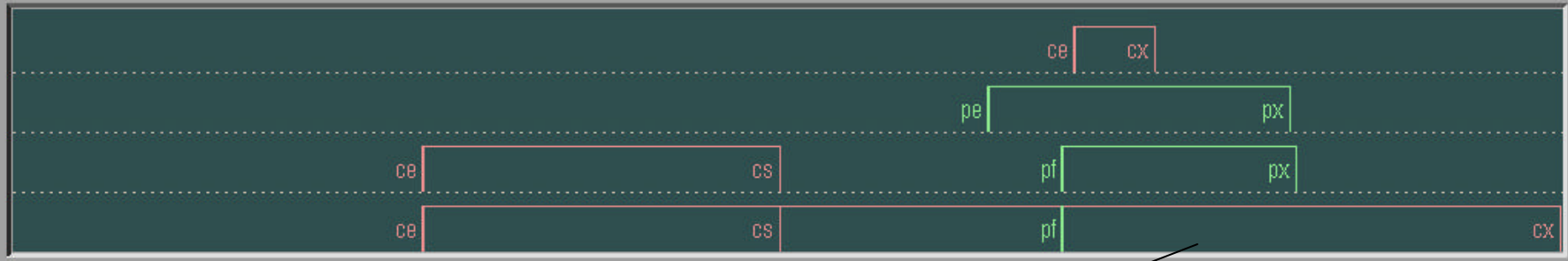
# Parse: Drop-Off

Action label

Component labels

Object tracks

```
Track-15: car_enter(790) SKIP car_exit(790) P = 0.04029918  
Event 42: 0.027689 5 42  
Segmentation:  
Person pass-through, frames 1998 - 2047  
P = 0.02008930  
Track-16: person_enter(780) SKIP person_exit(780) P = 0.02008931  
Event 43: 0.025000 10 43  
Segmentation:  
Person drove in, frames 1906 - 2048  
P = 0.02433817  
Track-17: car_enter(724) SKIP car_stop(724) SKIP person_found(787) SKIP person_exit(787) P = 0.01780757  
Event 45: 0.017689 13 45  
Segmentation:  
Person drop off, frames 1906 - 2097  
P = 0.01780757  
Track-18: car_enter(724) SKIP car_stop(724) SKIP person_found(787) SKIP car_exit(724) P = 0.0168
```



Temporal extent

Figures from: Ivanov, Yuri, Chris Stauffer, Aaron Bobick, and W. E. L. Grimson. "Video Surveillance of Interactions." *IEEE Workshop on Visual Surveillance (ICCV 2001)* (1999). Courtesy of IEEE, Y. Ivanov, C. Stauffer, A. Bobick, W. E. L. Grimson. Copyright 1999 IEEE. Used with Permission.