

Lecture 8

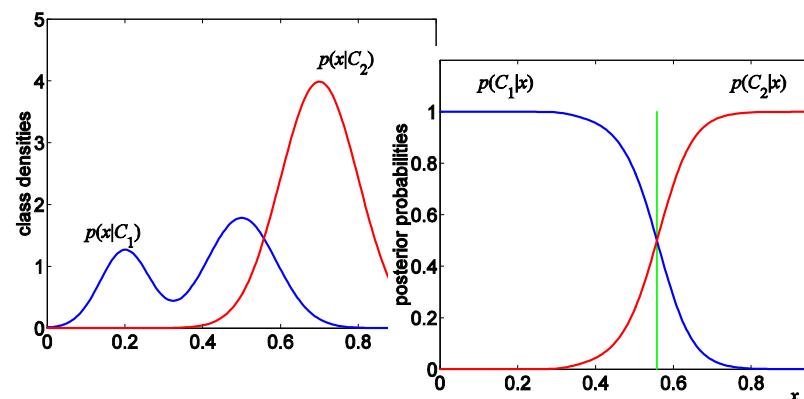
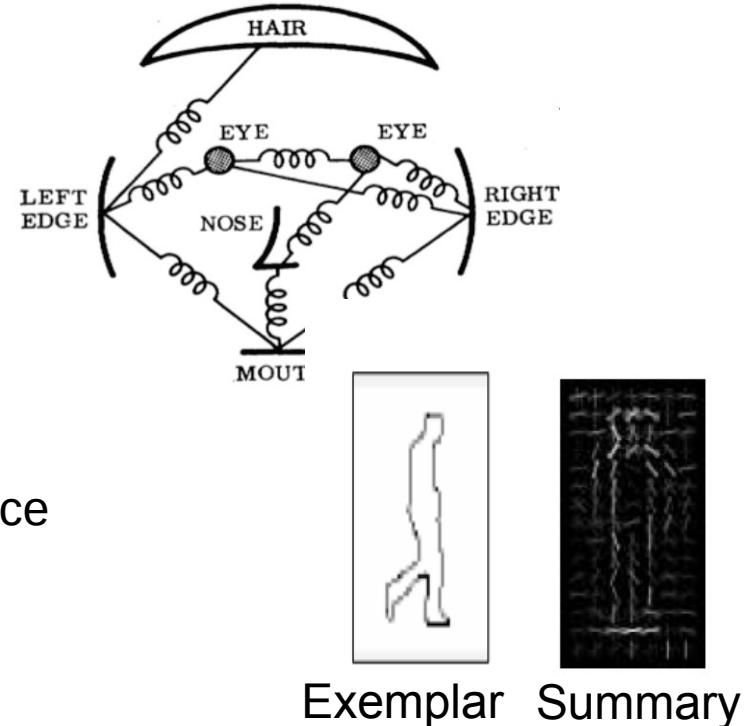
Statistical Learning & Part Models

COS 429: Computer Vision

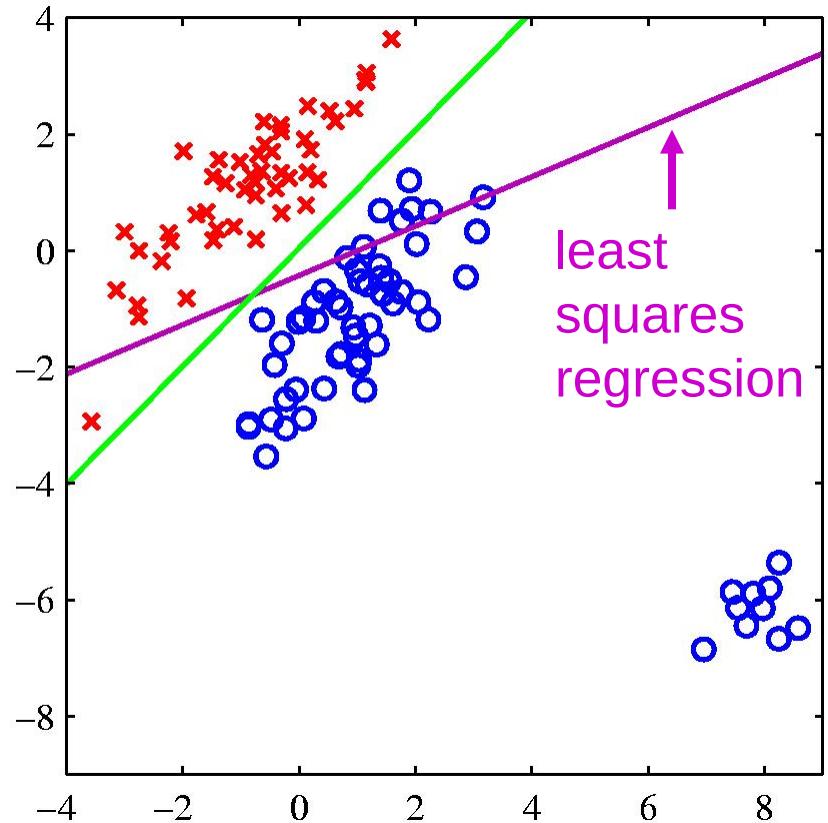
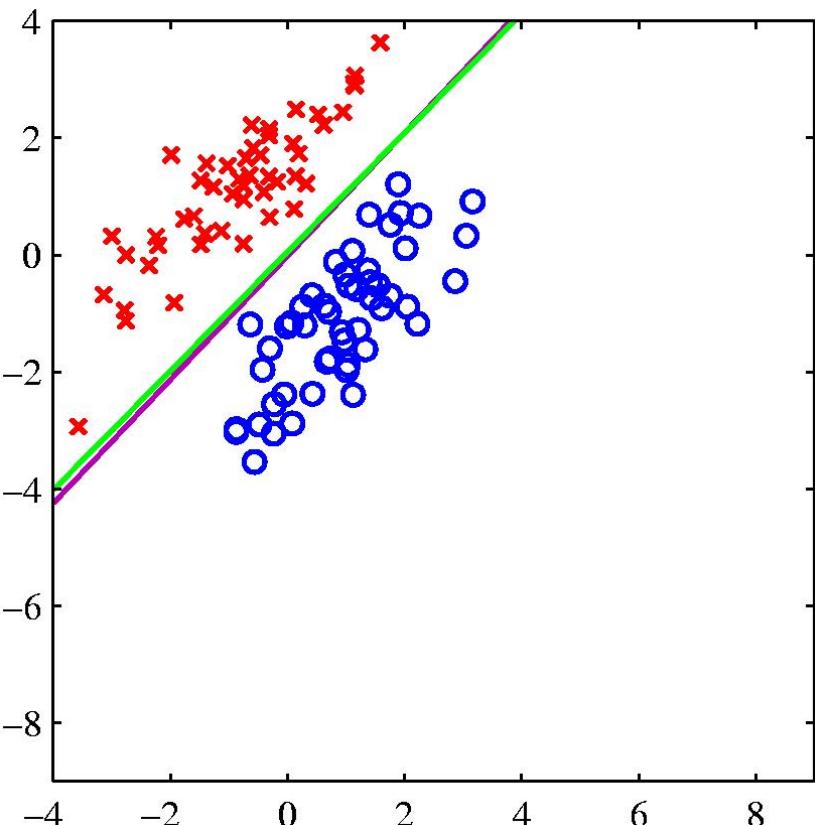


Review: Object Detection Typical Components

- **Hypothesis** generation
 - Sliding window, Segmentation, feature point detection, random, search
- **Encoding** of (local) image data
 - Colors, Edges, Corners, Histogram of Oriented Gradients, Wavelets, Convolution Filters
- **Relationship** of different parts to each other
 - Blur or histogram, Tree/Star, Pairwise/Covariance
- **Learning** from (labeled) examples
 - Selecting representative examples (templates), Clustering, Building a cascade
 - Classifiers: Bayes, Logistic regression, SVM, AdaBoost, Neural Network...
 - Generative vs. Discriminative
- **Verification** - removing redundant, overlapping, incompatible examples
 - Non-Max Suppression, context priors, geometry



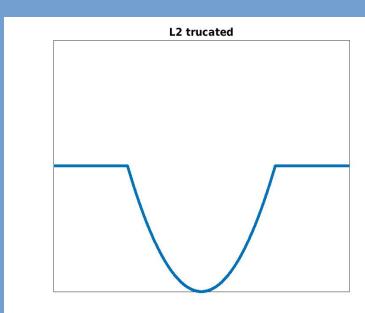
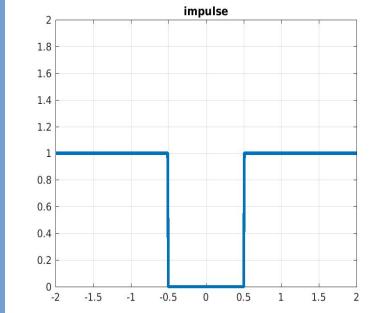
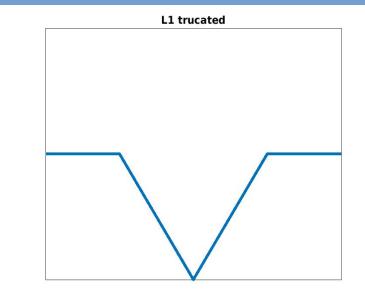
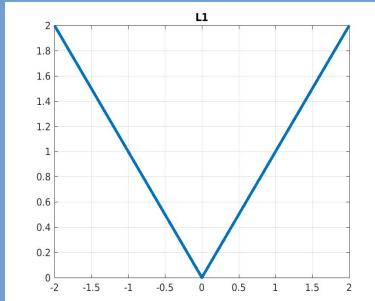
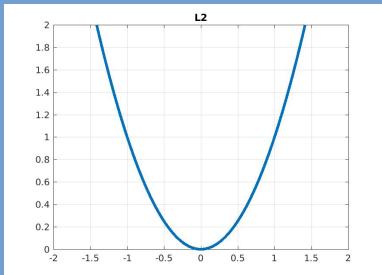
Using least squares for classification



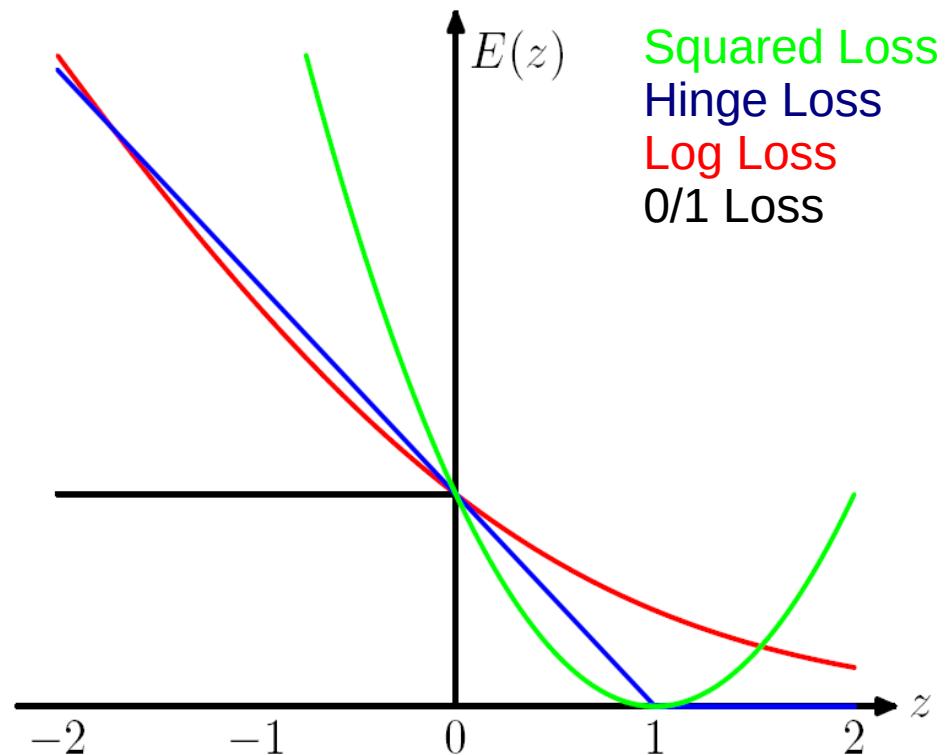
If the right answer is 1 and the model says 1.5, it loses, so it changes the boundary to avoid being “too correct”

The Problem: Loss Function

Recall: Regression Loss Functions



Some Classification Loss Functions



Sigmoid

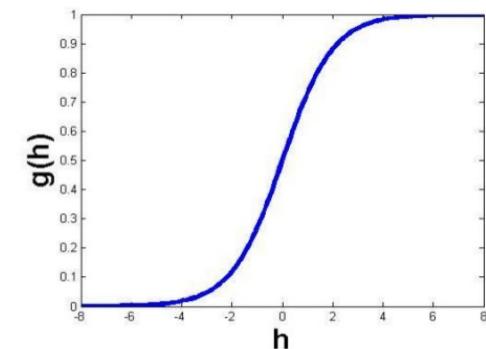
We model the probability of a label Y to be equal $y \in \{-1, 1\}$, given a data point $x \in \mathbf{R}^n$, as:

$$P(Y = y | x) = \frac{1}{1 + \exp(-y(w^T x + b))}.$$

This amounts to modeling the *log-odds ratio* as a linear function of X :

$$\log \frac{P(Y = 1 | x)}{P(Y = -1 | x)} = w^T x + b.$$

binary
 $\frac{1}{1 + e^{w^T x}}$
 $e^{w^T x}$



- ▶ The decision boundary $P(Y = 1 | x) = P(Y = -1 | x)$ is the hyperplane with equation $w^T x + b = 0$.
- ▶ The region $P(Y = 1 | x) \geq P(Y = -1 | x)$ (i.e., $w^T x + b \geq 0$) corresponds to points with predicted label $\hat{y} = +1$.

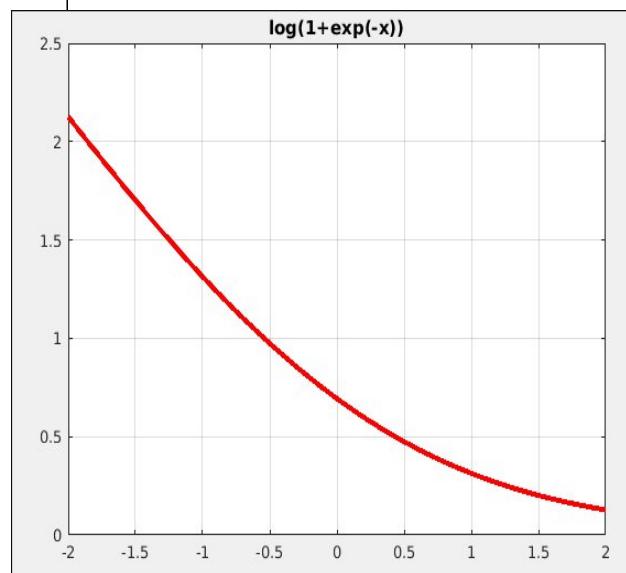
Log Loss

The likelihood function is

$$l(w, b) = \prod_{i=1}^m \frac{1}{1 + e^{-y_i(w^T x_i + b)}}.$$

Now maximize the log-likelihood:

$$\max_{w,b} L(w, b) := - \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i + b)})$$



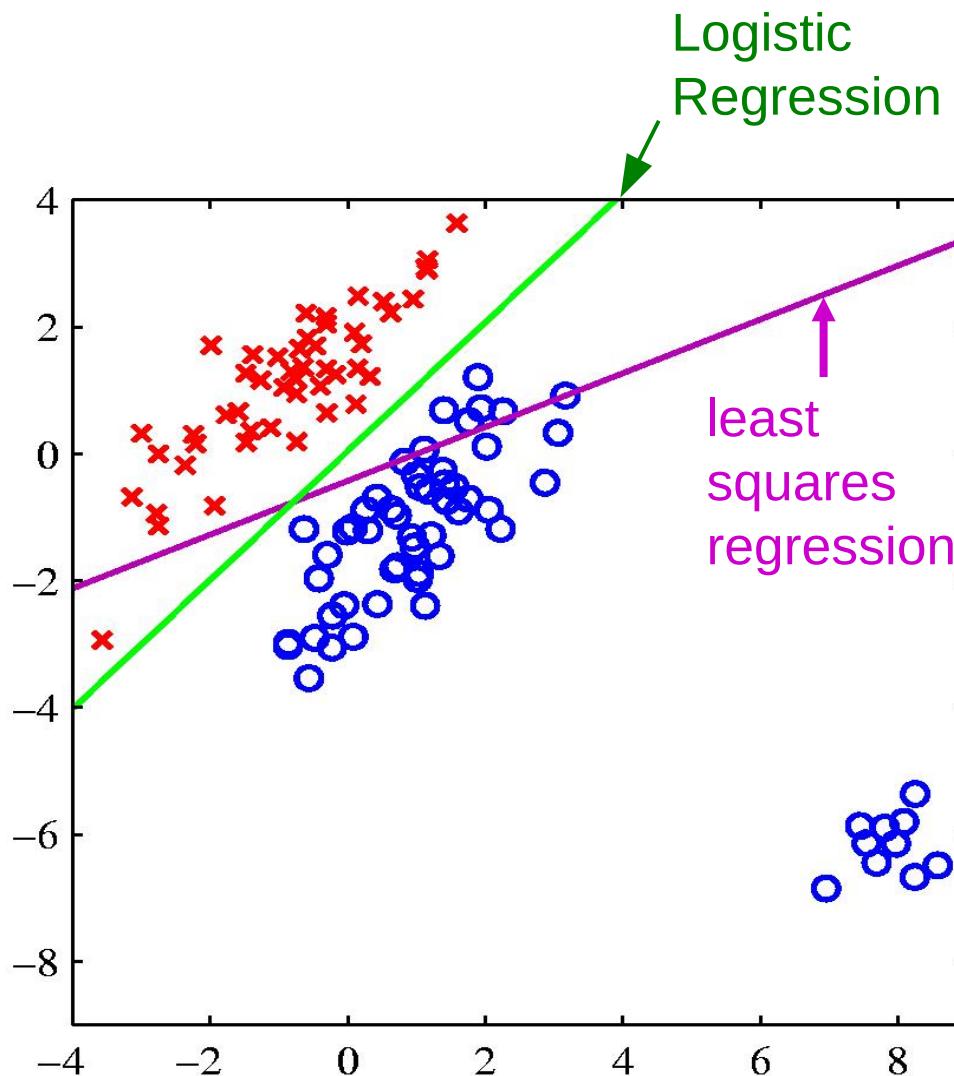
In practice, we may consider adding a regularization term

$$\max_{w,b} L(w, b) + \lambda r(w),$$

with $r(w) = \|w\|_2^2$ or $r(w) = \|w\|_1$.

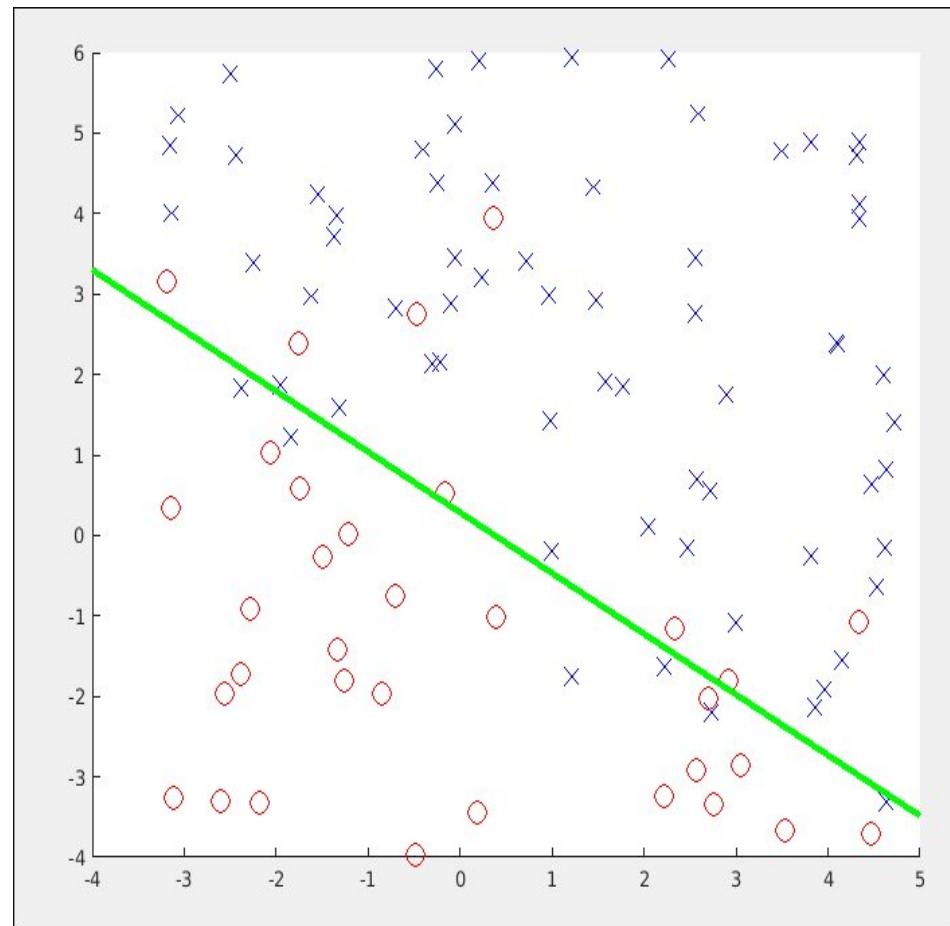
The Logistic Loss:
 $\log(1 + \exp(-z))$
Where
 $z = y^*(w^T x + b)$

Logistic Result

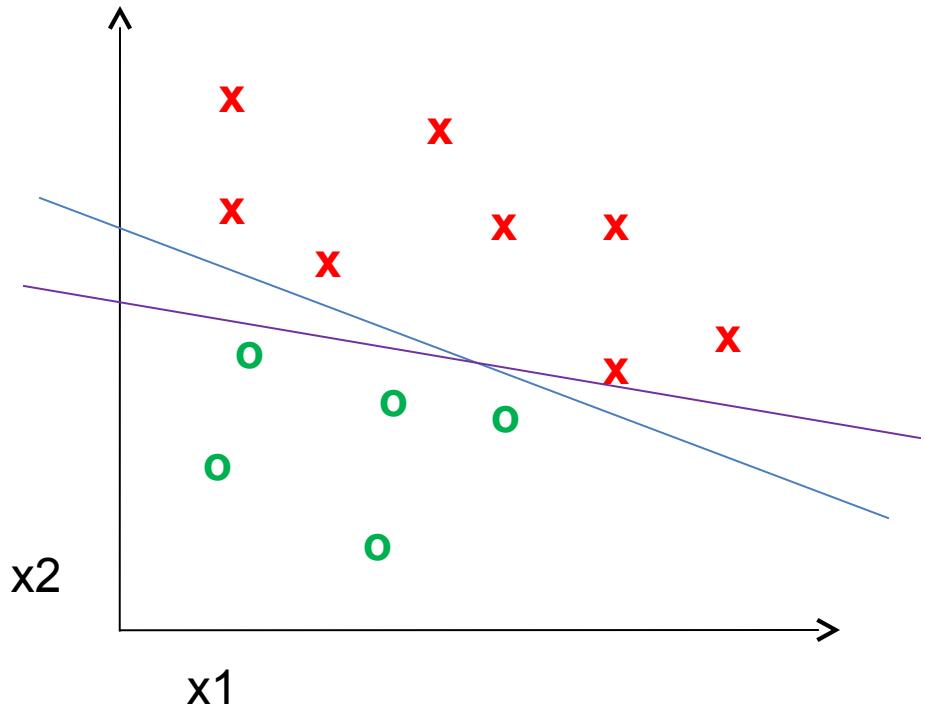


Using Logistic Regression

- Quick, simple classifier (try it first)
- Outputs a probabilistic label confidence
- Use L2 or L1 regularization
 - L1 does feature selection and is robust to irrelevant features but slower to train



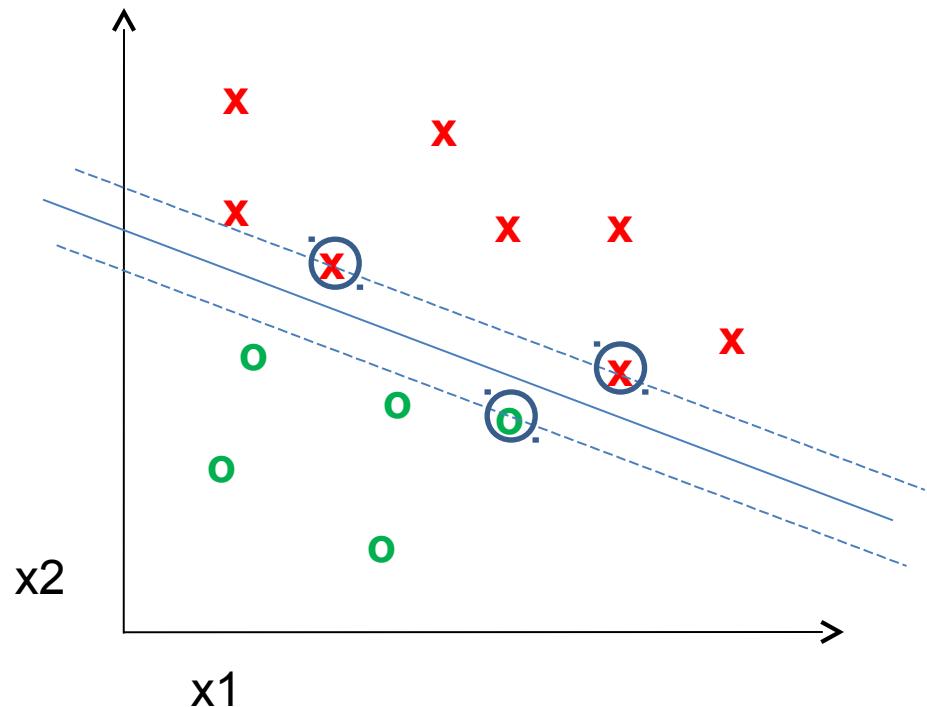
Classifiers: Linear SVM



- Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

Classifiers: Linear SVM

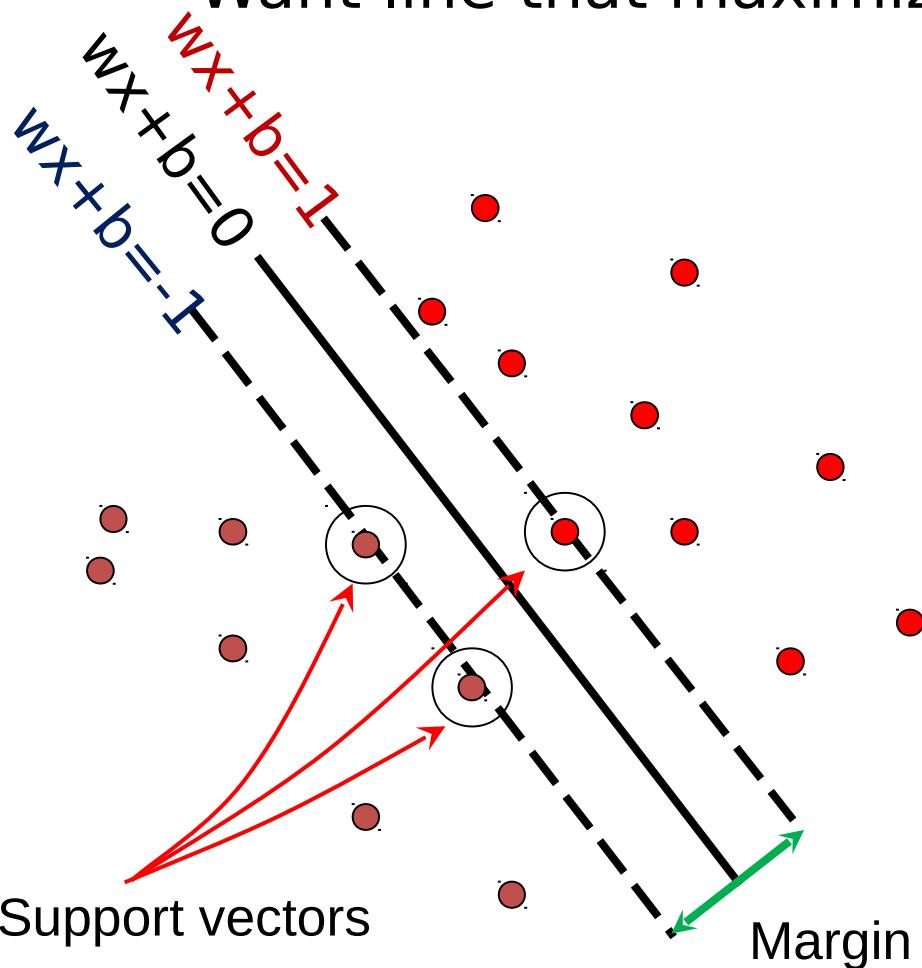


- Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

Support vector machines: Margin

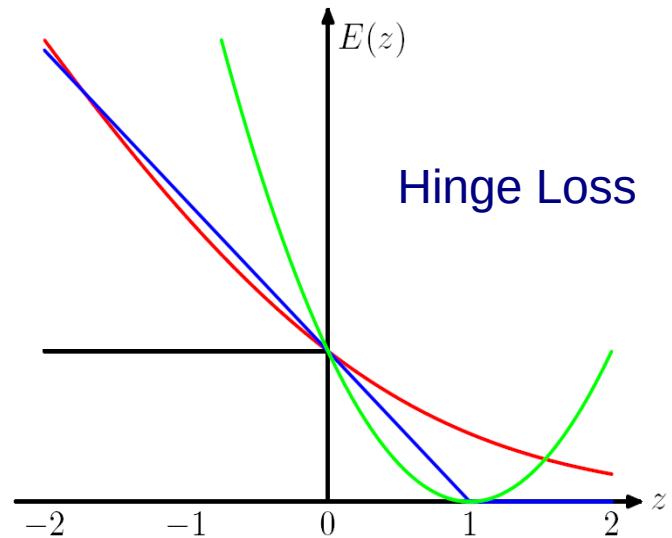
- Want line that maximizes the margin.



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

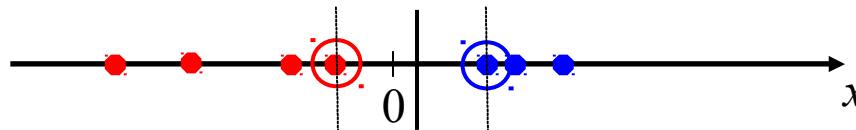
For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$



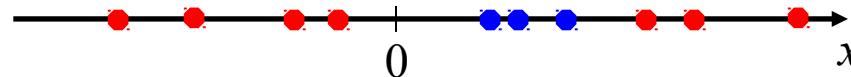
$$L(y, f(x)) = \max(0, 1 - y \cdot f(x))$$

Nonlinear SVMs

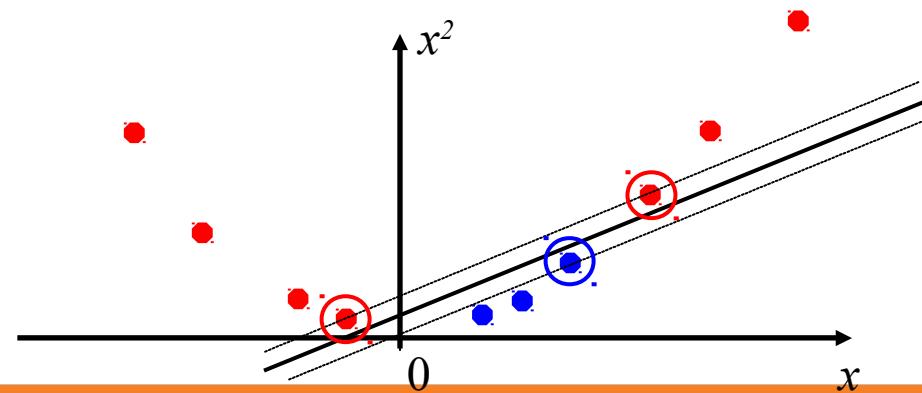
- Datasets that are linearly separable work out great:



- But what if the dataset is just too hard?

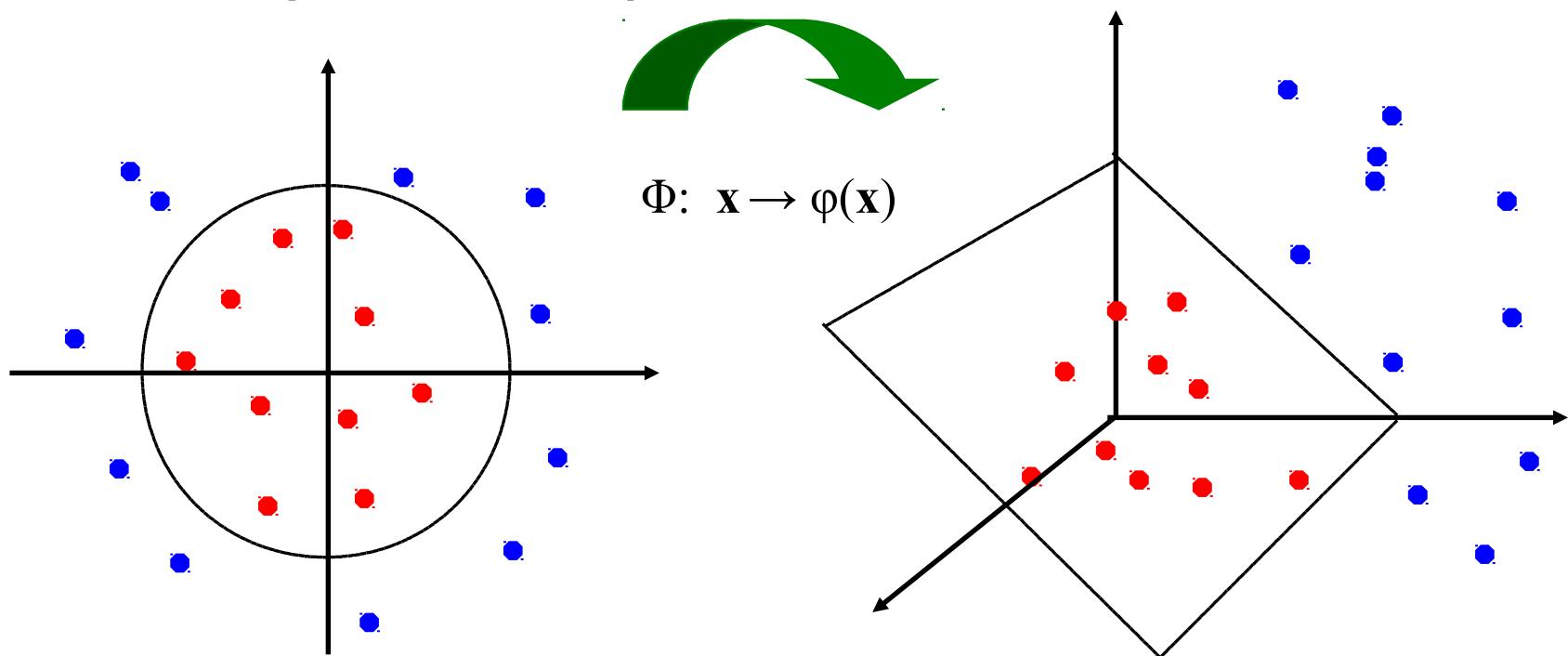


- We can map it to a higher-dimensional space:



Nonlinear SVMs

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



Nonlinear SVMs

- *The kernel trick:* instead of explicitly computing the lifting transformation $\phi(\mathbf{x})$, define a kernel function K such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

(to be valid, the kernel function must satisfy *Mercer's condition*)

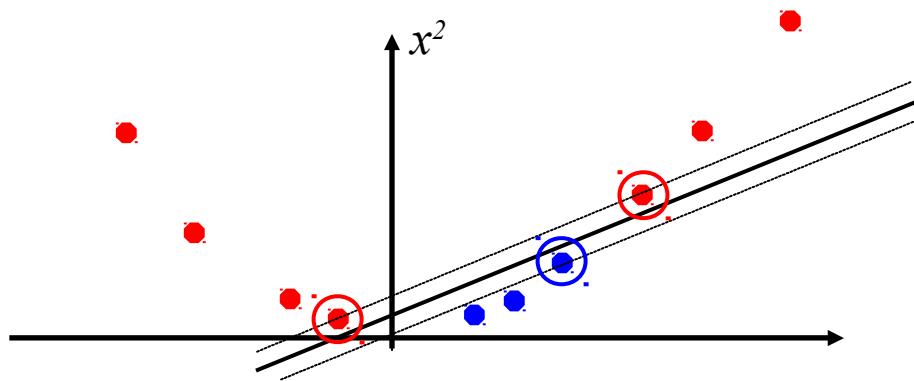
- This gives a nonlinear decision boundary in the original feature space:

$$\sum_i \alpha_i y_i \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

C. Burges, [A Tutorial on Support Vector Machines for Pattern Recognition](#), Data Mining and Knowledge Discovery, 1998

Nonlinear kernel: Example

- Consider the mapping $\phi(x) = (x, x^2)$



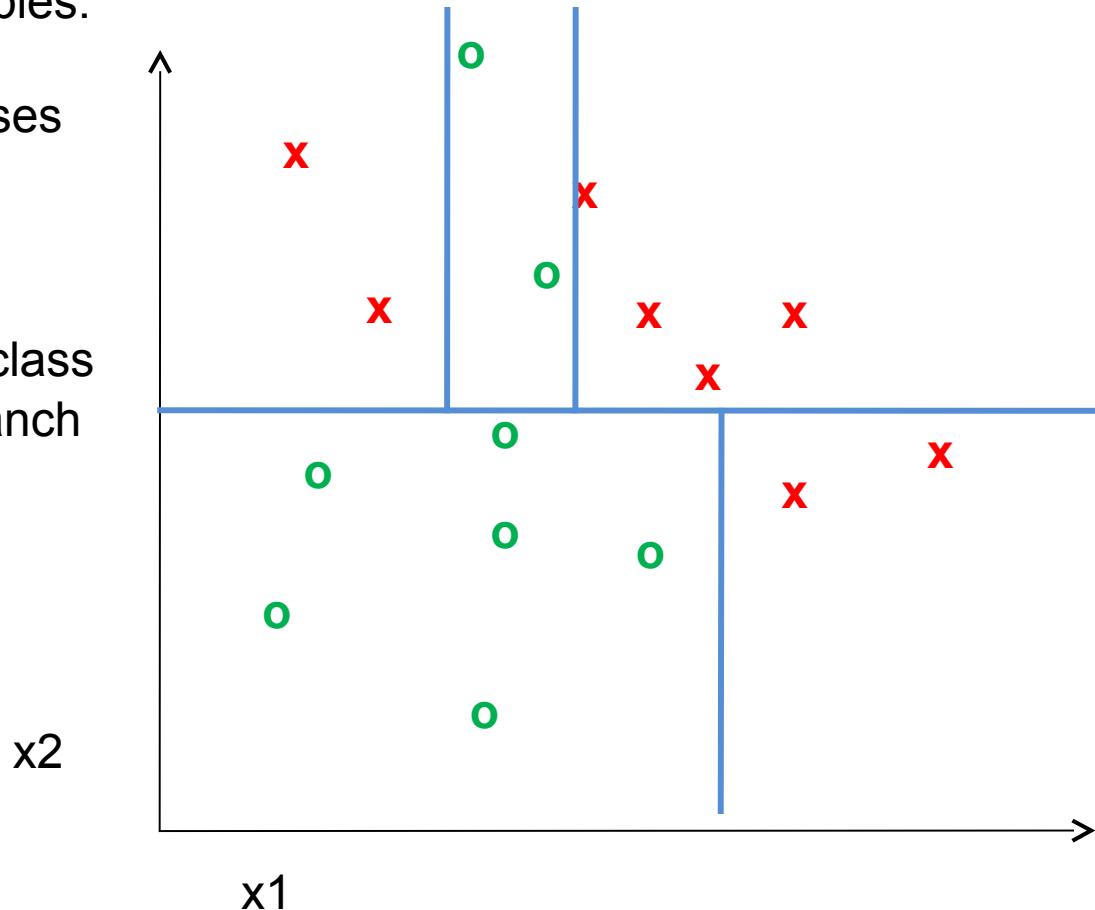
$$\phi(x) \cdot \phi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2y^2$$

$$K(x, y) = xy + x^2y^2$$

Classifiers: Decision Trees

Given (weighted) labeled examples:

- Select best single feature & threshold that separates classes
- For each branch, recurse
- Stop when
 - some depth is reached
 - Branch is (close to) single-class
 - too few examples left in branch

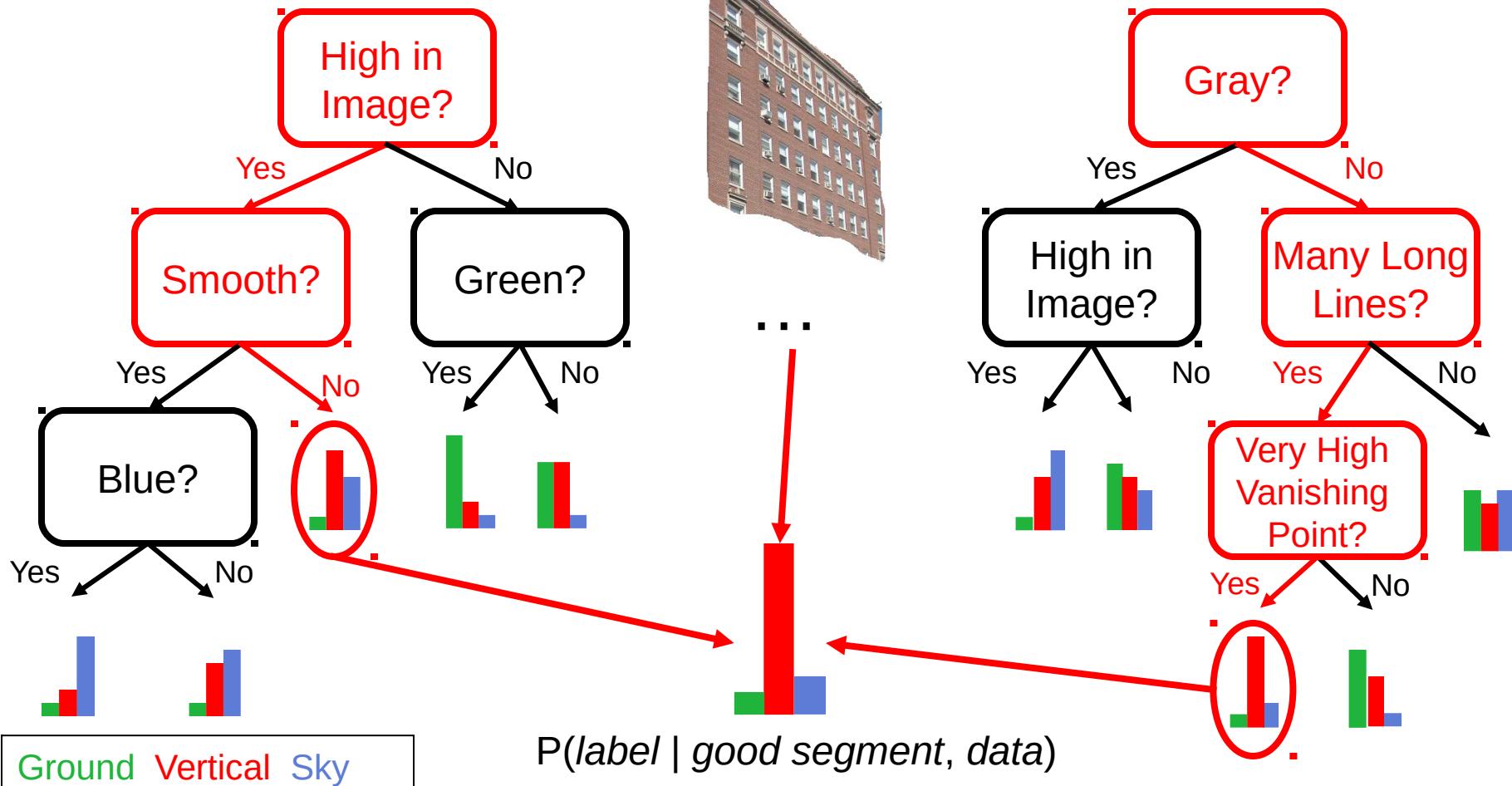


Ensemble Methods: Boosting

Discrete AdaBoost(Freund & Schapire 1996b)

1. Start with weights $w_i = 1/N$, $i = 1, \dots, N$.
2. Repeat for $m = 1, 2, \dots, M$:
 - (a) Fit the classifier $f_m(x) \in \{-1, 1\}$ using weights w_i on the training data.
 - (b) Compute $\text{err}_m = E_w[1_{(y \neq f_m(x))}]$, $c_m = \log((1 - \text{err}_m)/\text{err}_m)$.
 - (c) Set $w_i \leftarrow w_i \exp[c_m \cdot 1_{(y_i \neq f_m(x_i))}]$, $i = 1, 2, \dots, N$, and renormalize so that $\sum_i w_i = 1$.
3. Output the classifier $\text{sign}[\sum_{m=1}^M c_m f_m(x)]$

Boosted Decision Trees

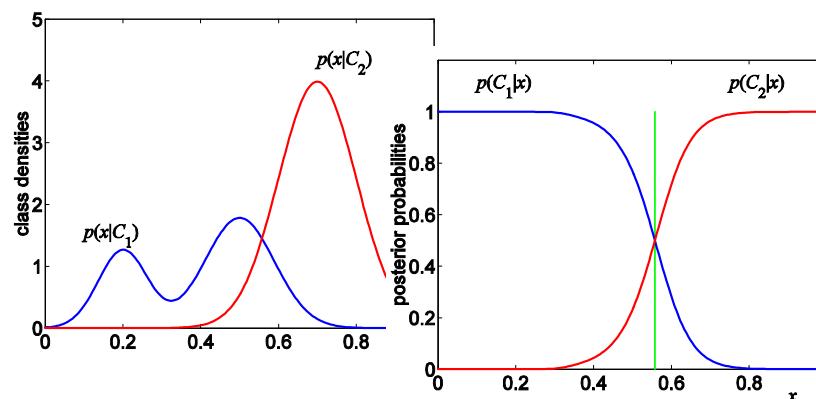
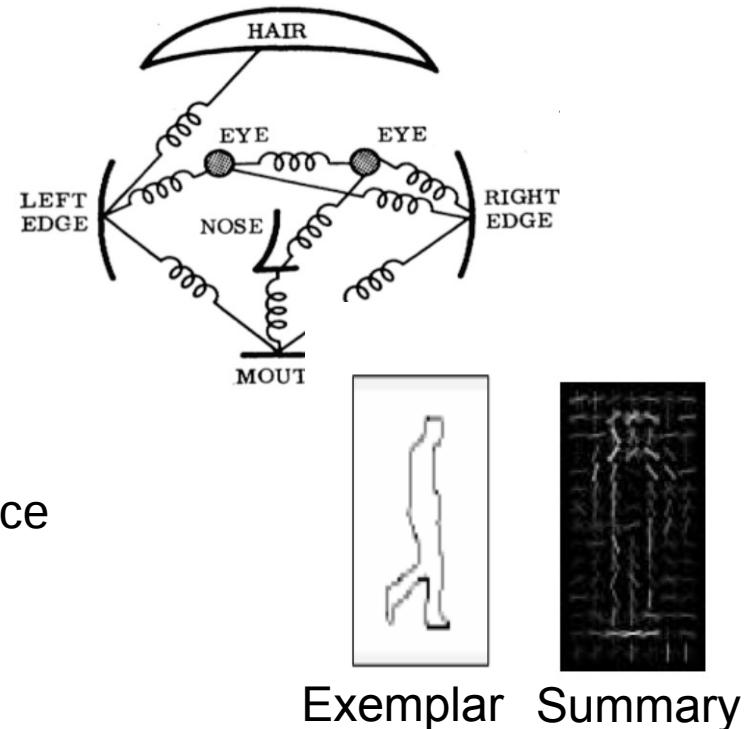


Using Boosted Decision Trees

- Flexible: can deal with both continuous and categorical variables
- How to control bias/variance trade-off
 - Size of trees
 - Number of trees
- Boosting trees often works best with a small number of well-designed features
- Boosting “stubs” can give a fast classifier

Review: Typical Components

- **Hypothesis** generation
 - Sliding window, Segmentation, feature point detection, random, search
- **Encoding** of (local) image data
 - Colors, Edges, Corners, Histogram of Oriented Gradients, Wavelets, Convolution Filters
- **Relationship** of different parts to each other
 - Blur or histogram, Tree/Star, Pairwise/Covariance
- **Learning** from labeled examples
 - Selecting representative examples (templates), Clustering, Building a cascade
 - Classifiers: Bayes, Logistic regression, SVM, Decision Trees, AdaBoost, ...
 - **Generative vs. Discriminative**
- **Verification** - removing redundant, overlapping, incompatible examples
 - Non-Max Suppression, context priors, geometry



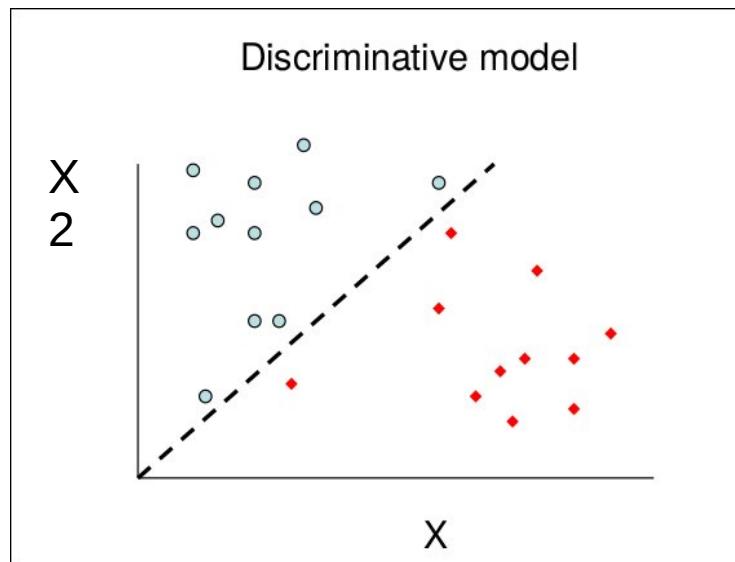
Discriminative vs. Generative Classifiers

Training classifiers involves estimating $f: X \rightarrow Y$, or $P(Y|X)$ “Y given X”

Discriminative Classification:

Find the boundary between classes

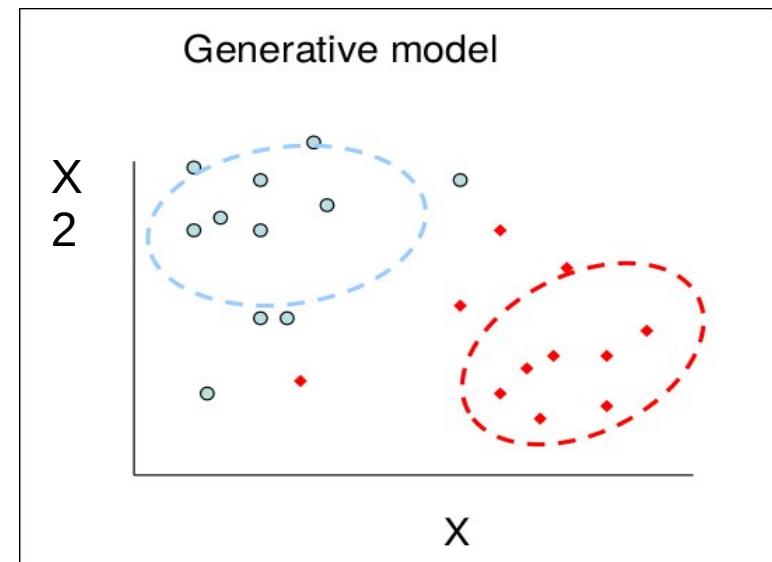
1. Assume some functional form for $P(Y|X)$
2. Estimate parameters of $P(Y|X)$ directly from training data



Generative Classification:

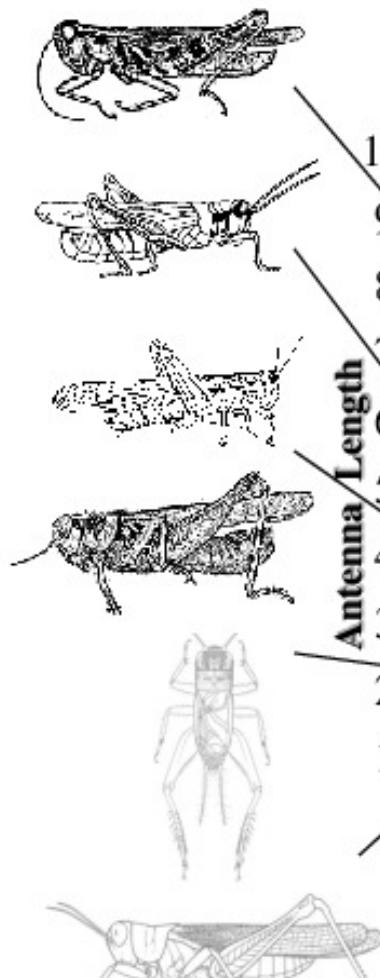
Model each class & see which fits better

1. Assume some functional form for $P(X|Y)$, $P(X)$
2. Estimate parameters of $P(X|Y)$, $P(X)$ directly from training data
3. Use Bayes rule to calculate $P(Y|X=x_i)$

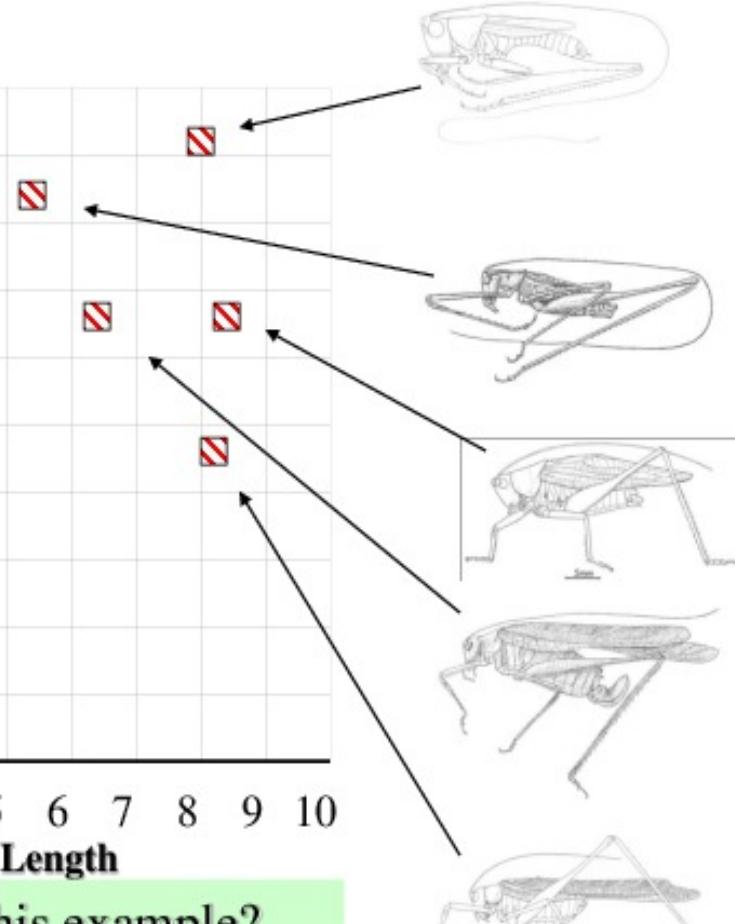


Bayesian Classification (Generative Model)

Grasshoppers



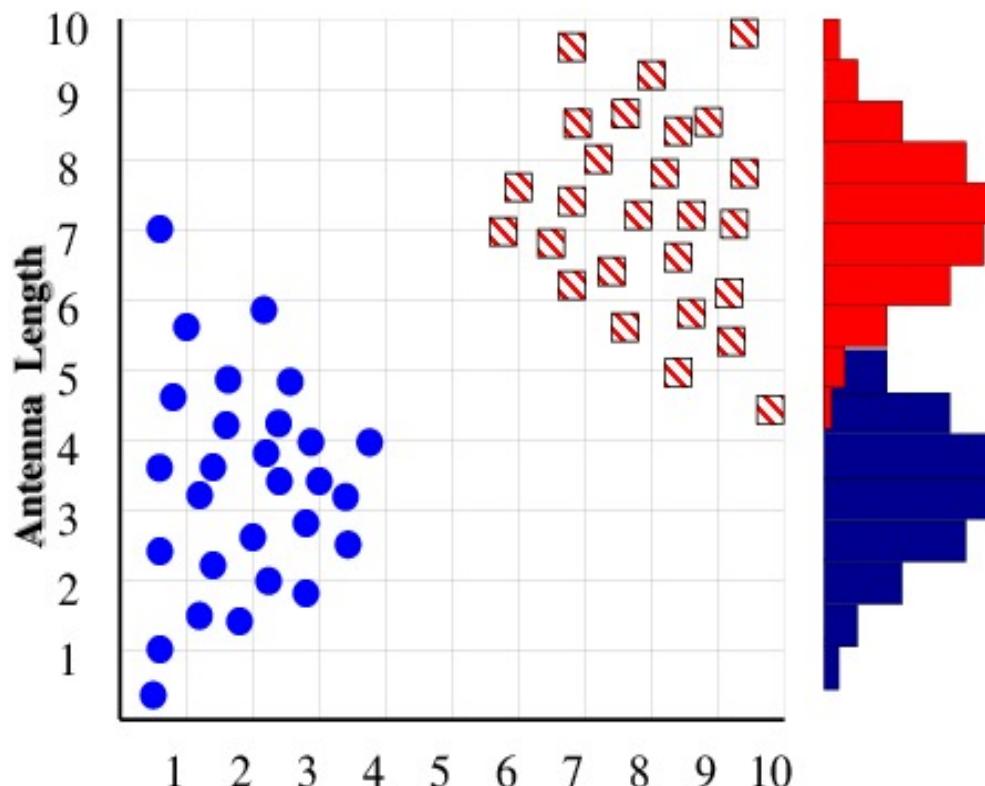
Katydid



Remember this example?

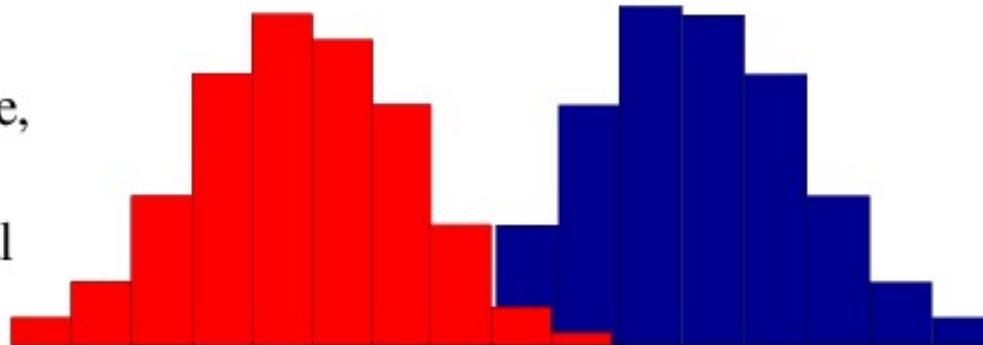
Bayesian Classification

With a lot of data, we can build a histogram. Let us just build one for “Antenna Length” for now....

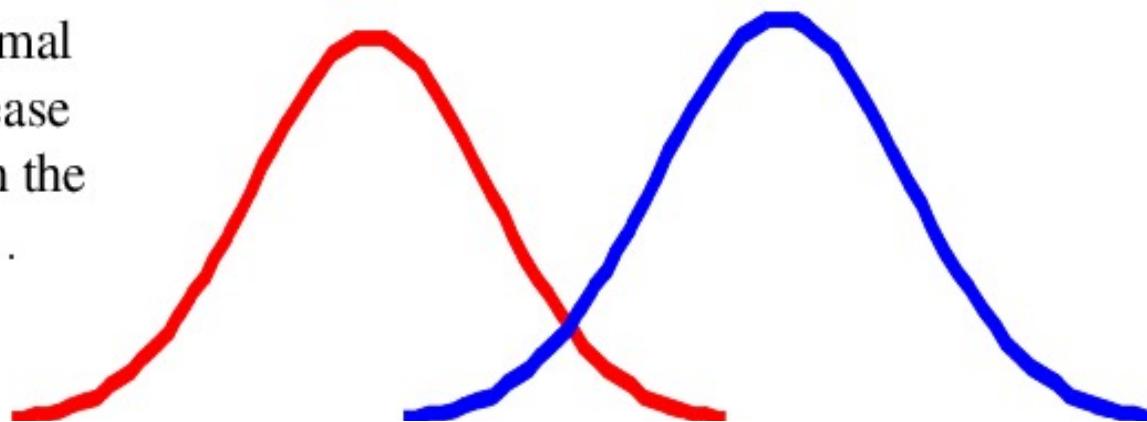


Bayesian Classification

We can leave the histograms as they are, or we can summarize them with two normal distributions.



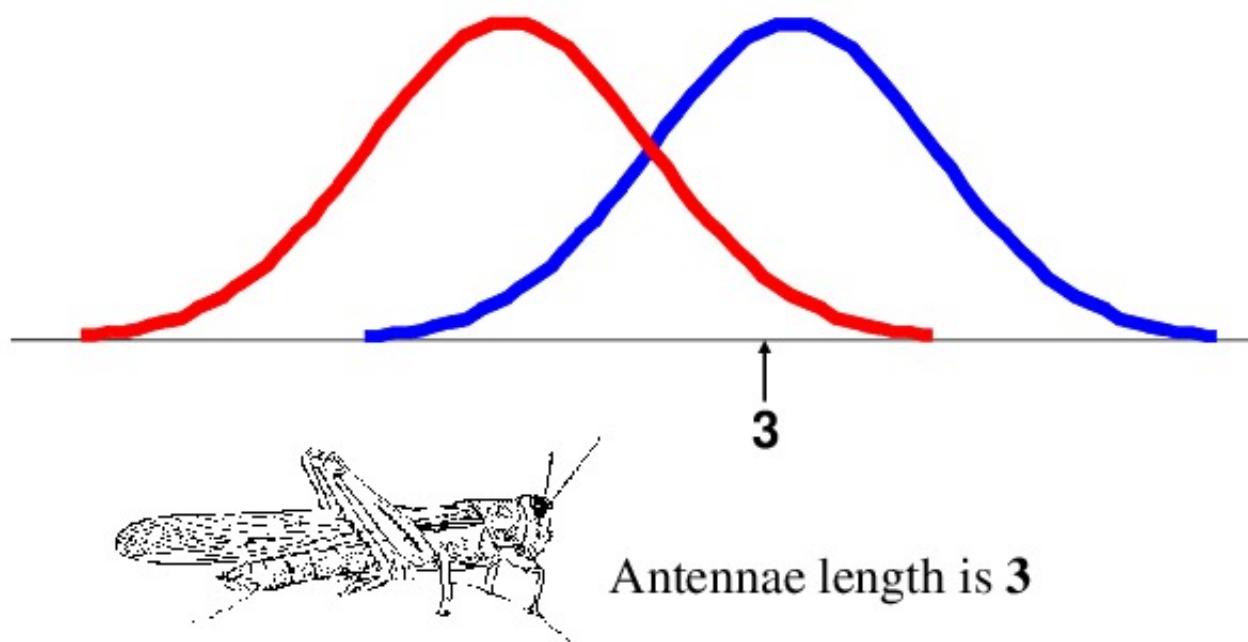
Let us use two normal distributions for ease of visualization in the following slides...



Bayesian Classification

- We want to classify an insect we have found. Its antennae are 3 units long. How can we classify it?
- We can just ask ourselves, give the distributions of antennae lengths we have seen, is it more probable that our insect is a **Grasshopper** or a **Katydid**.
- There is a formal way to discuss the most probable classification...

$p(c_j | d)$ = probability of class c_j , given that we have observed d

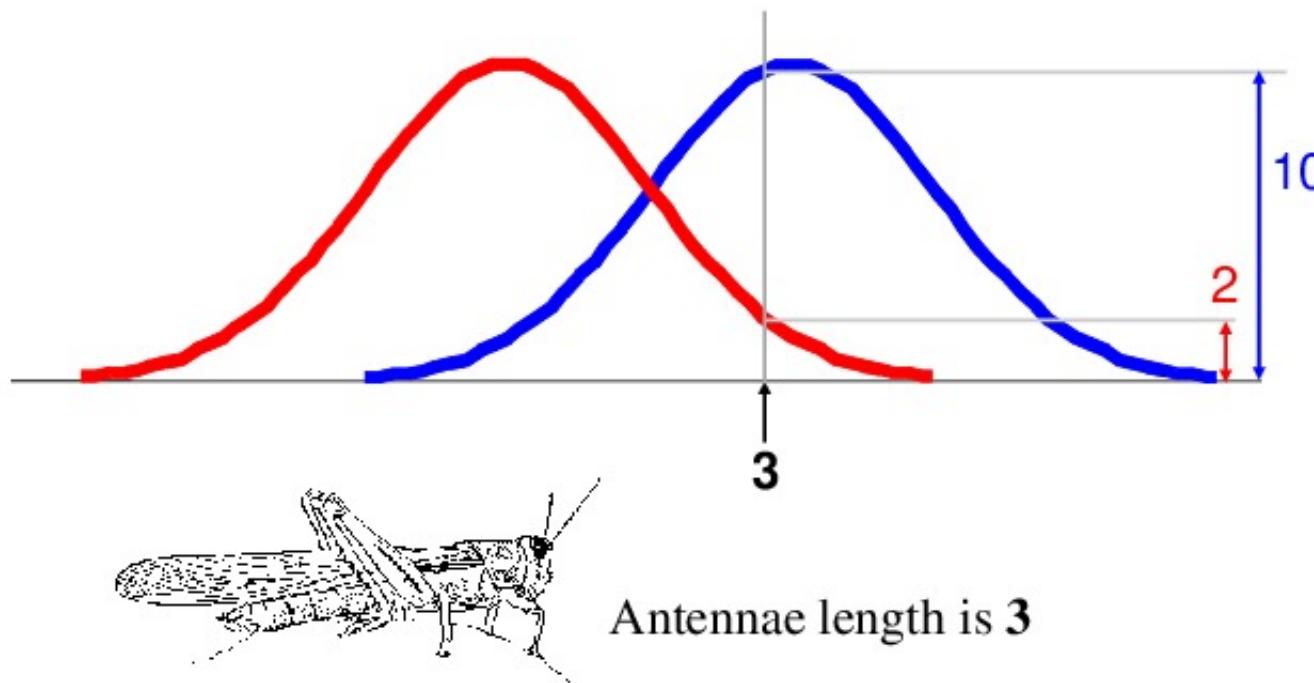


Bayesian Classification

$p(c_j | d)$ = probability of class c_j , given that we have observed d

$$P(\text{Grasshopper} | 3) = 10 / (10 + 2) = 0.833$$

$$P(\text{Katydid} | 3) = 2 / (10 + 2) = 0.166$$



Bayesian Classification

- Bayesian classifiers use **Bayes theorem**, which says

$$p(c_j | d) = \frac{p(d | c_j) p(c_j)}{p(d)}$$

- $p(c_j | d)$ = probability of instance d being in class c_j ,

This is what we are trying to compute

- $p(d | c_j)$ = probability of generating instance d given class c_j ,

We can imagine that being in class c_j , causes you to have feature d with some probability

- $p(c_j)$ = probability of occurrence of class c_j ,

This is just how frequent the class c_j , is in our database

- $p(d)$ = probability of instance d occurring

This can actually be ignored, since it is the same for all classes

Naïve Bayes Classification

- To simplify the task, **naïve Bayesian classifiers** assume attributes have independent distributions, and thereby estimate

$$p(d|c_j) = p(d_1|c_j) * p(d_2|c_j) * \dots * p(d_n|c_j)$$



The probability of class c_j generating instance d , equals....



The probability of class c_j generating the observed value for feature 1, multiplied by..



The probability of class c_j generating the observed value for feature 2, multiplied by..



Naïve Bayes, Odds ratio, and Logit

Assuming independence of features d_1 and d_2 , we can classify between 2 classes $\{C_1, C_2\}$, by compute the ratio:

$$\frac{P(C_1|d_1, d_2)}{P(C_2|d_1, d_2)} = \frac{P(d_1|C_1) P(d_2|C_1) P(C_1)}{P(d_1|C_2) P(d_2|C_2) P(C_2)} < 1$$

This is called the Odds ratio.

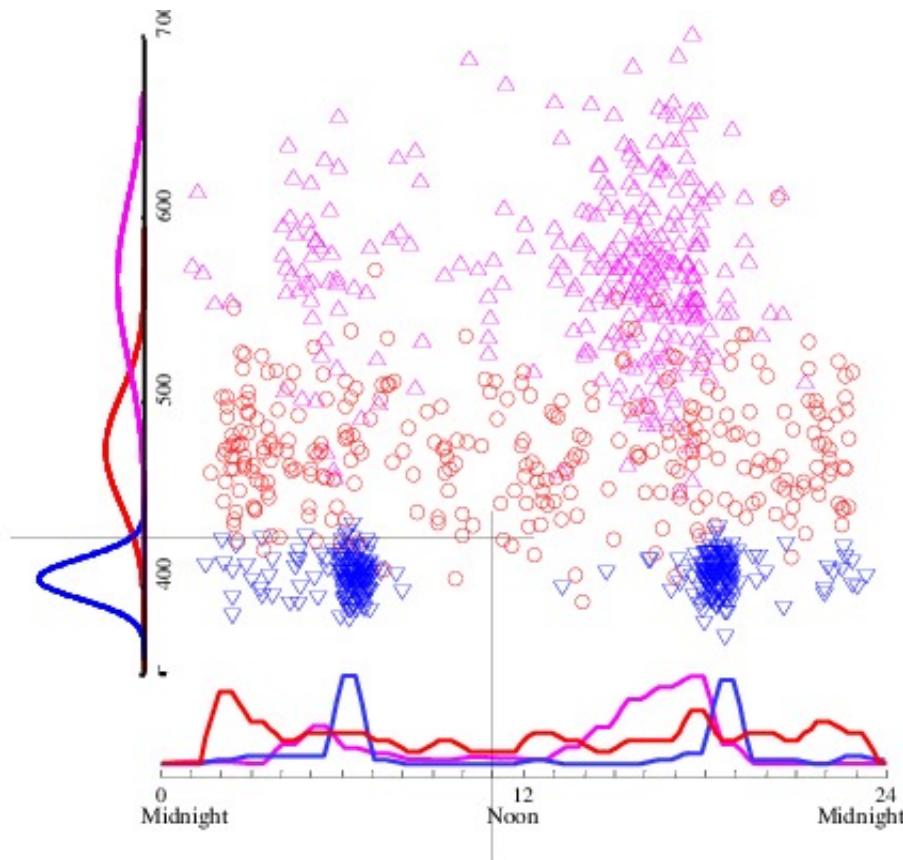
It is often easier to take the log of this, called the Log Odds or Logit:

$$\begin{aligned} \log \frac{P(C_1|d_1, d_2)}{P(C_2|d_1, d_2)} \\ = \log \frac{P(d_1|C_1)}{P(d_1|C_2)} + \log \frac{P(d_2|C_1)}{P(d_2|C_2)} + \log \frac{P(C_1)}{P(C_2)} < 0 \end{aligned}$$

Naïve Bayes

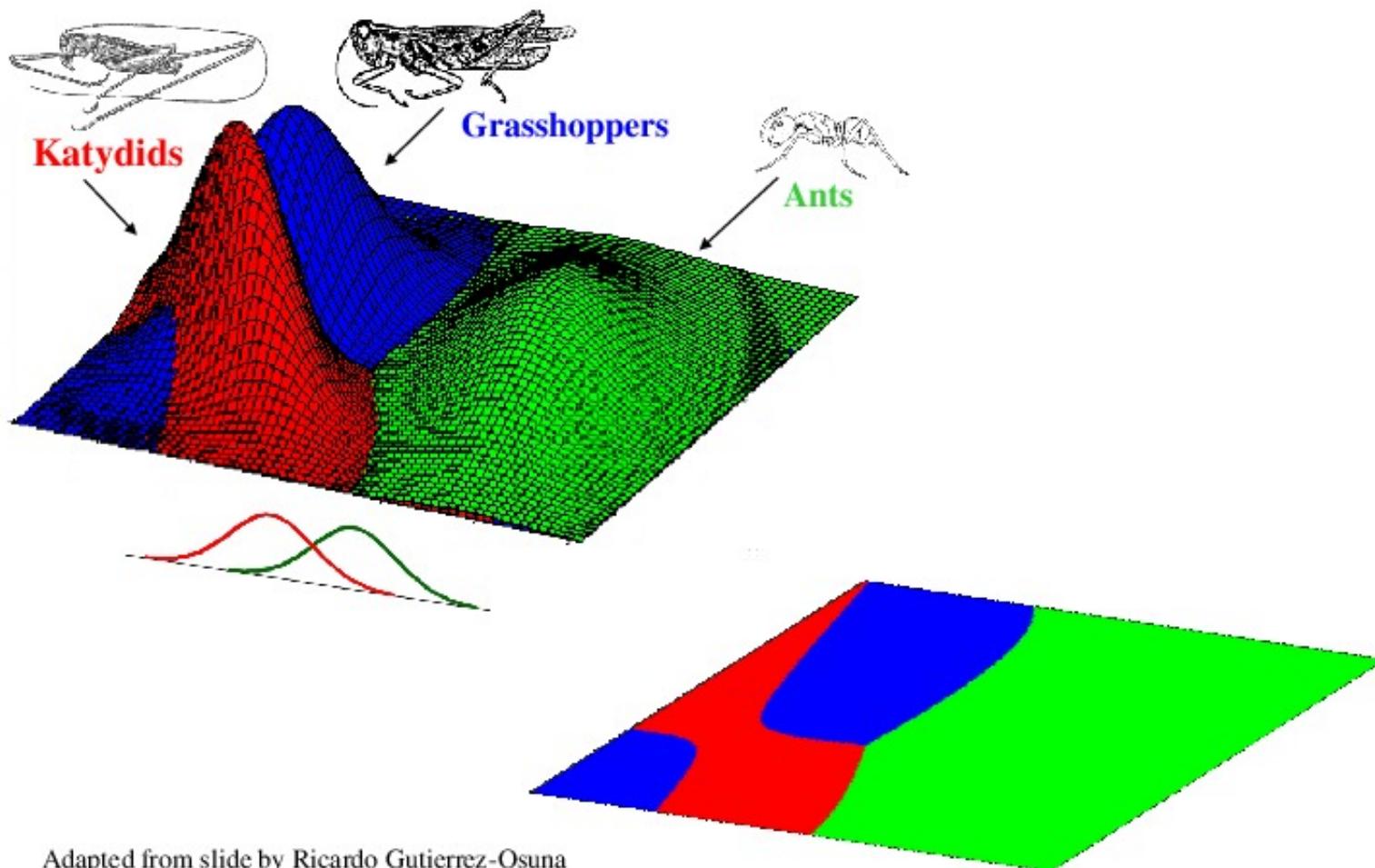
Suppose I observe an insect with a wingbeat frequency of 420 at 11:00am

What is it?



Naïve Bayes

Naive Bayes with Gaussian densities have piecewise quadratic decision boundary.



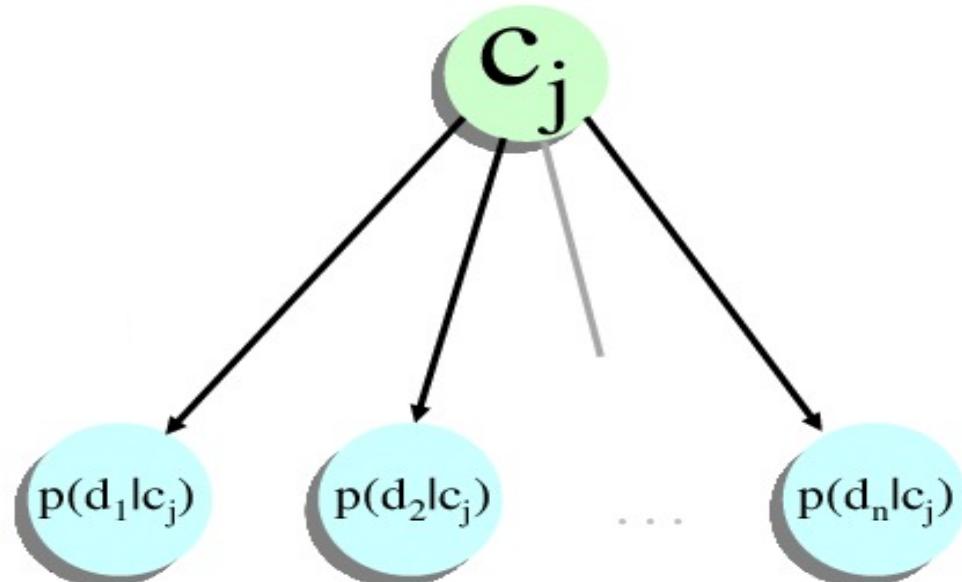
Adapted from slide by Ricardo Gutierrez-Osuna

Graphical Models

Independence assumption of Naive Bayes is the simplest assumption. Often much more complicated relationships needs to be represented.

A good way to do this is a Graphical Model
(aka. Bayesian Network)

Naive Bayes assumes independence of d_1, d_2, \dots conditioned on the class c



Animal	Mass > 10 kg	
Cat	Yes	0.15
	No	0.85
Dog	Yes	0.91
	No	0.09
Pig	Yes	0.99
	No	0.01

Animal	Color	
Cat	Black	0.33
	White	0.23
	Brown	0.44
Dog	Black	0.97
	White	0.03
	Brown	0.90
Pig	Black	0.04
	White	0.01

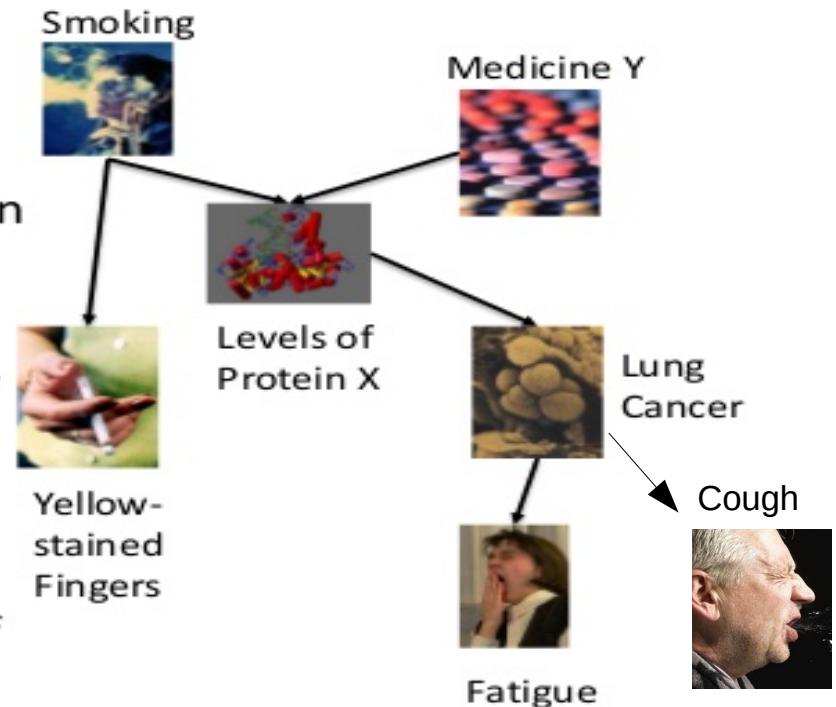
Animal
Cat
Dog
Pig

Graphical Models

More complicated models consider the dependence between different variables.

Bayesian Network

1. Factorize the jpa
2. Answer questions like:
 1. $P(\text{Lung Cancer} | \text{Levels of Protein } X) = ?$
 2. $\text{Ind}(\text{Smoking}, \text{Fatigue} | \text{Levels of Protein } X) ?$
 3. What will happen if I design a drug that blocks the function of protein X (**predict effect of interventions**)?



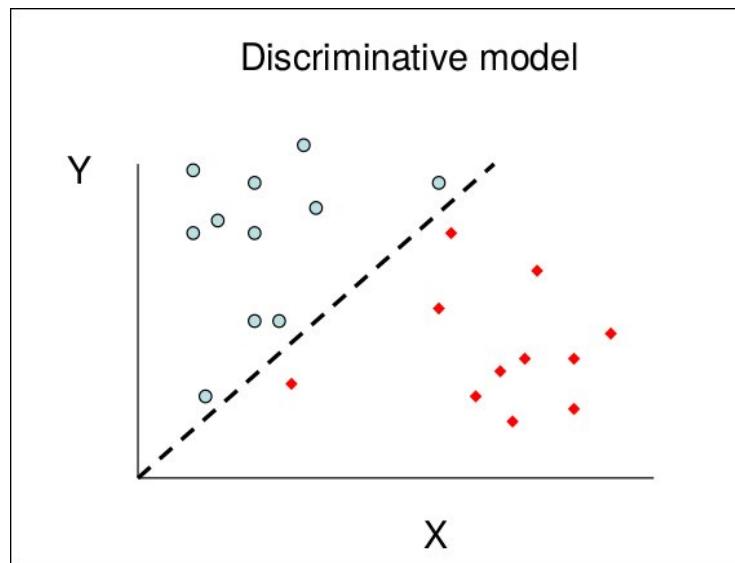
Discriminative vs. Generative Classifiers

Training classifiers involves estimating $f: X \rightarrow Y$, or $P(Y|X)$ “Y given X”

Discriminative Classification:

Find the boundary between classes

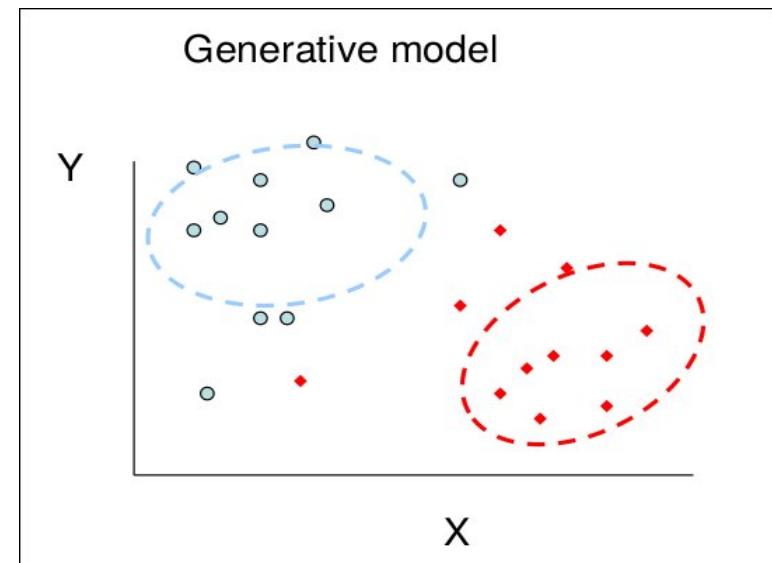
1. Assume some functional form for $P(Y|X)$
2. Estimate parameters of $P(Y|X)$ directly from training data



Generative Classification:

Model each class & see which fits better

1. Assume some functional form for $P(X|Y)$, $P(X)$
2. Estimate parameters of $P(X|Y)$, $P(X)$ directly from training data
3. Use Bayes rule to calculate $P(Y|X=x_i)$



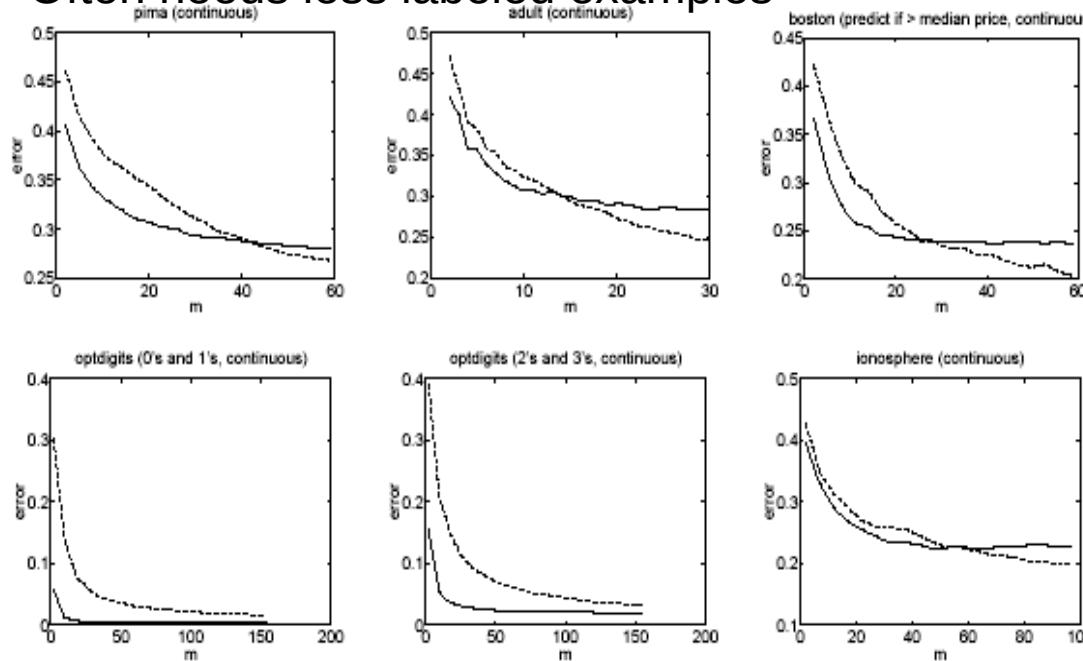
Discriminative vs. Generative Classifiers

Discriminative:

- + Model directly what you care about
- + With many examples usually more accurate
- + Often faster to evaluate, can scale well to many examples & classes

Generative:

- + Allows more flexibility to model relationships between variables
- + Can handle compositionality, missing & occluded parts (more “object oriented”)
- + Often needs less labeled examples



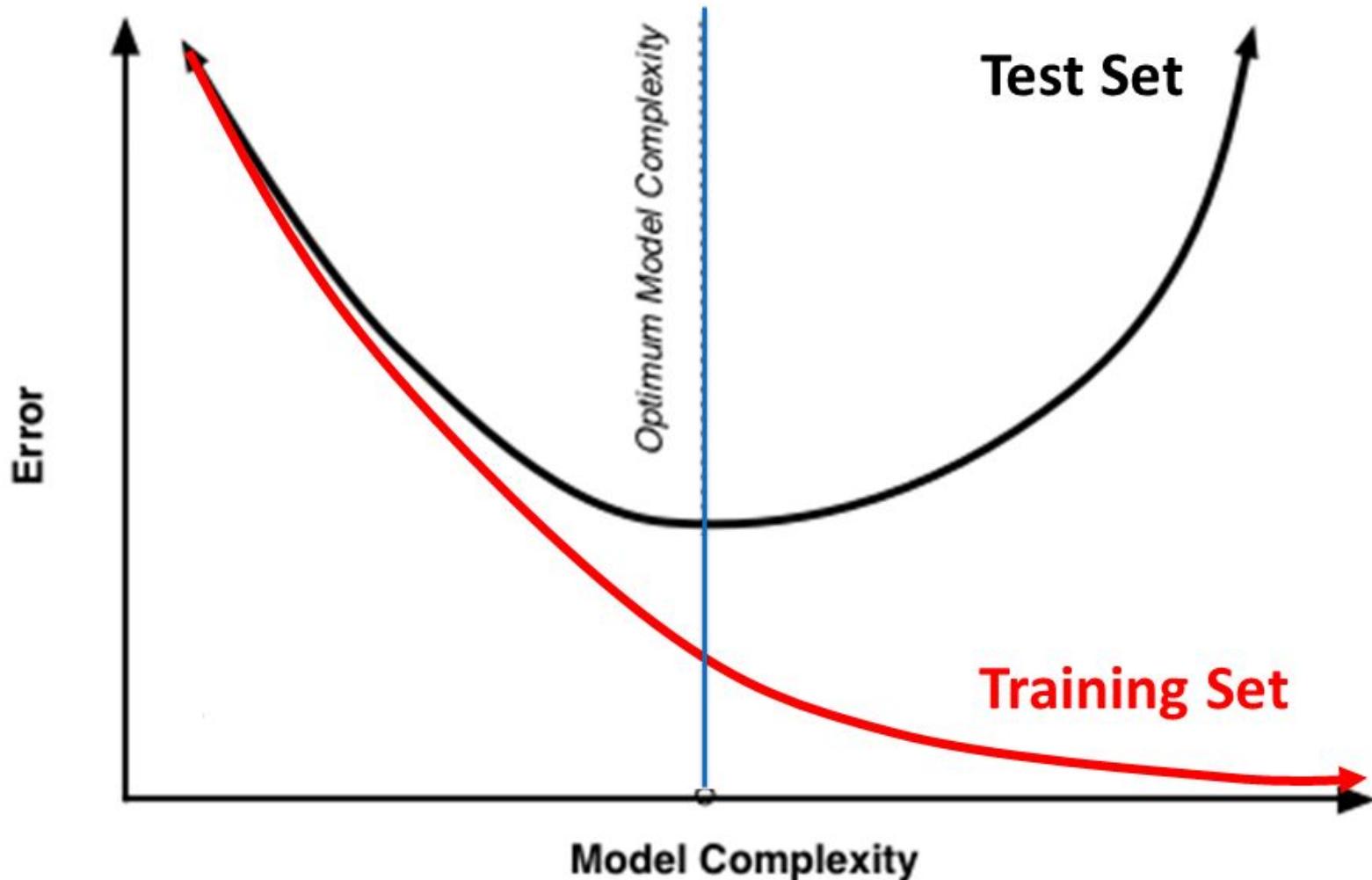
“On Discriminative vs. Generative classifiers: A comparison of logistic regression and naïve Bayes,” A. Ng and M. Jordan, NIPS 2002.

Comparison

assuming \mathbf{x} in $\{0, 1\}$

	Learning Objective	Training	Inference
Naïve Bayes	$\text{maximize} \sum_i \left[\sum_j \log P(x_{ij} y_i; \theta_j) \right] + \log P(y_i; \theta_0)$	$\theta_{kj} = \frac{\sum_i \delta(x_{ij} = 1 \wedge y_i = k) + r}{\sum_i \delta(y_i = k) + Kr}$	$\theta_1^T \mathbf{x} + \theta_0^T (1 - \mathbf{x}) > 0$ where $\theta_{1j} = \log \frac{P(x_{1j} = 1 y=1)}{P(x_{1j} = 1 y=0)}$, $\theta_{0j} = \log \frac{P(x_{0j} = 0 y=1)}{P(x_{0j} = 0 y=0)}$
Logistic Regression	$\text{maximize} \sum_i \log(P(y_i \mathbf{x}, \boldsymbol{\theta})) + \lambda \ \boldsymbol{\theta}\ $ where $P(y_i \mathbf{x}, \boldsymbol{\theta}) = 1/(1 + \exp(-y_i \boldsymbol{\theta}^T \mathbf{x}))$	Gradient ascent	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Linear SVM	$\text{minimize} \lambda \sum_i \xi_i + \frac{1}{2} \ \boldsymbol{\theta}\ ^2$ such that $y_i \boldsymbol{\theta}^T \mathbf{x} \geq 1 - \xi_i \quad \forall i$	Linear programming	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Kernelized SVM	complicated to write	Quadratic programming	$\sum_i y_i \alpha_i K(\hat{\mathbf{x}}_i, \mathbf{x}) > 0$
Nearest Neighbor	most similar features \rightarrow same label	Record data	y_i where $i = \operatorname{argmin}_i K(\hat{\mathbf{x}}_i, \mathbf{x})$

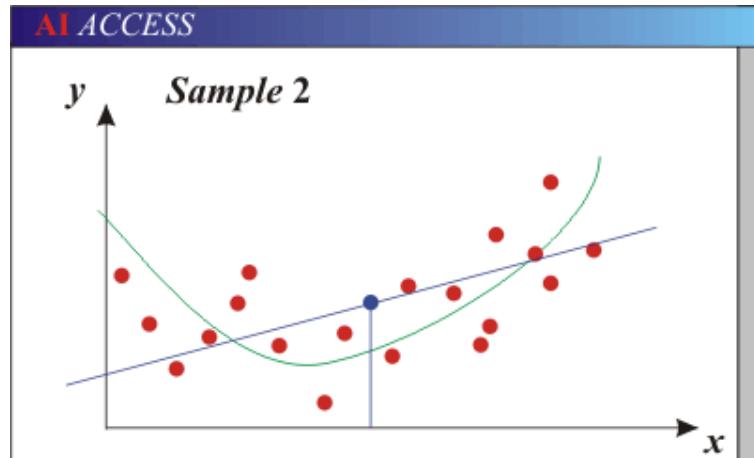
Train vs. Test Accuracy



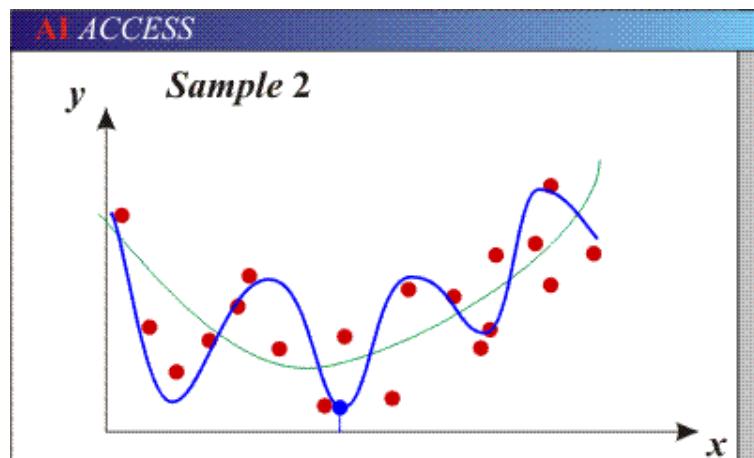
Generalization

- Components of generalization error
 - **Bias:** how much the average model over all training sets differ from the true model?
 - Error due to inaccurate assumptions/simplifications made by the model
 - **Variance:** how much models estimated from different training sets differ from each other
- **Underfitting:** model is too “simple” to represent all the relevant class characteristics
 - High bias and low variance
 - High training error and high test error
- **Overfitting:** model is too “complex” and fits irrelevant characteristics (noise) in the data
 - Low bias and high variance
 - Low training error and high test error

Bias-Variance Trade-off



- Models with too few parameters are inaccurate because of a large bias (not enough flexibility).
- Models with too many parameters are inaccurate because of a large variance (too much sensitivity to the sample).



Bias-Variance Trade-off

$$E(MSE) = \text{noise}^2 + \text{bias}^2 + \text{variance}$$

Unavoidable
error

Error due to
incorrect
assumptions

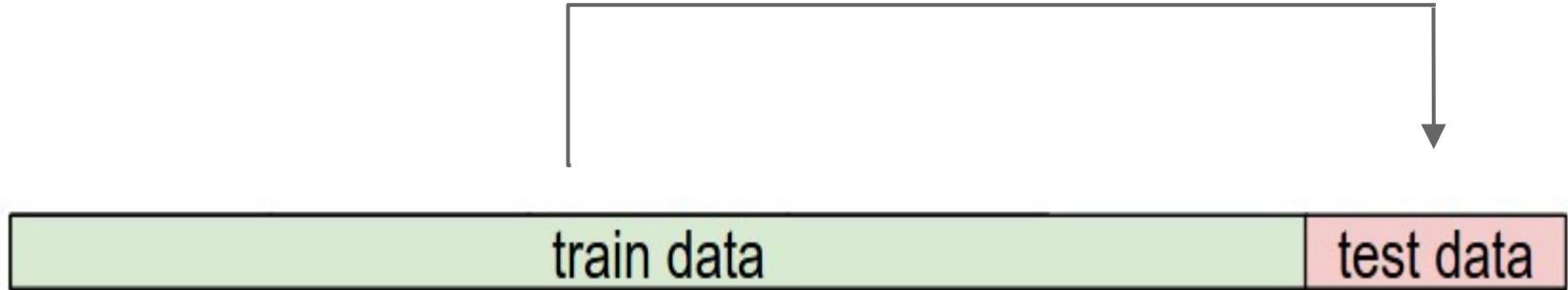
Error due to
variance of training
samples

See the following for explanations of bias-variance (also Bishop's "Neural Networks" book):

- <http://www.inf.ed.ac.uk/teaching/courses/mlsc/Notes/Lecture4/BiasVariance.pdf>

Cross Validation

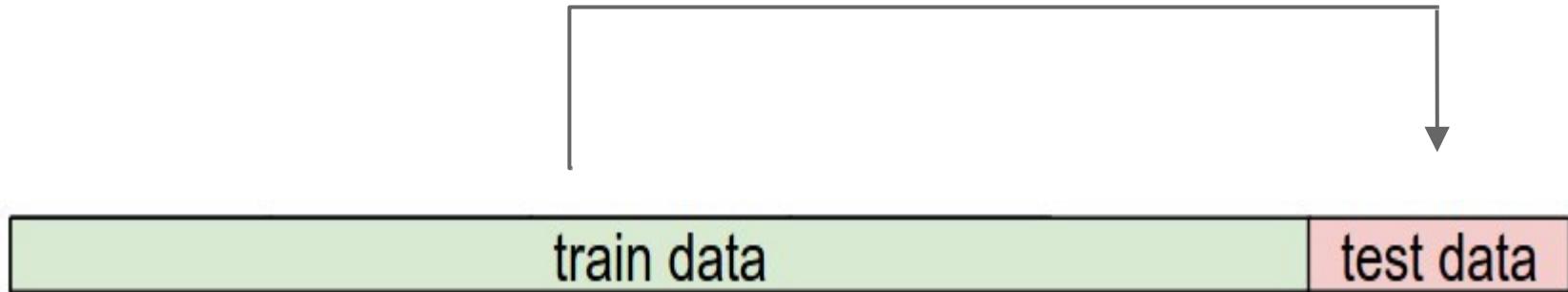
Try out what hyperparameters work best on test set.



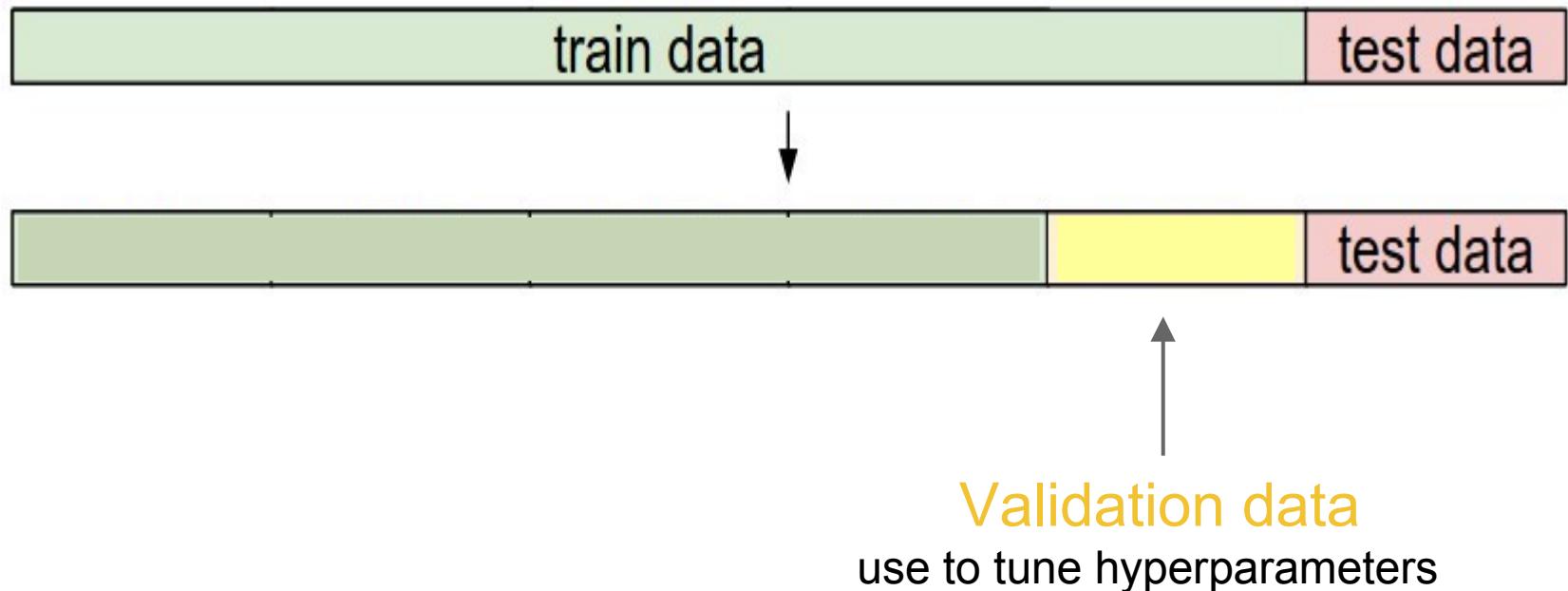
Cross Validation

Trying out what hyperparameters work best on test set:

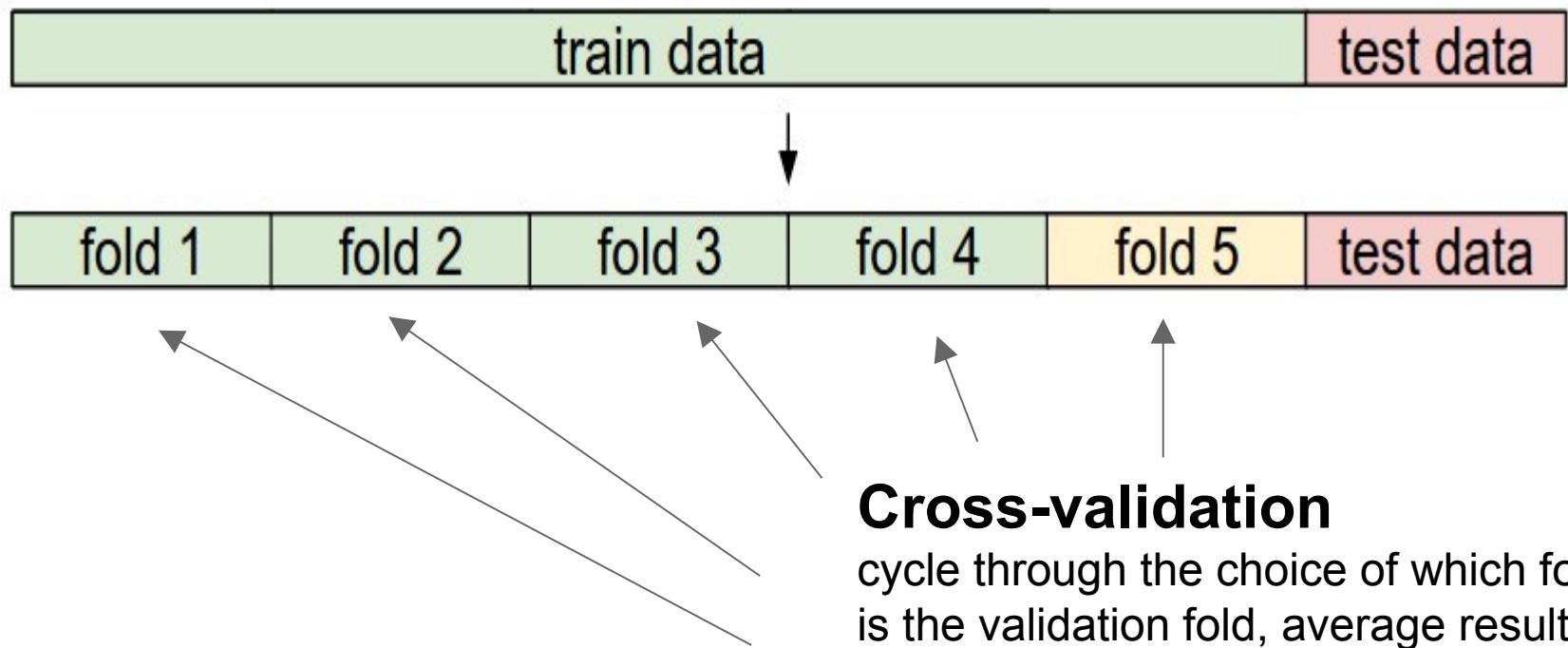
Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY**, at the end.



Validation



Cross Validation



So...

- No classifier is inherently better than any other: you need to make assumptions to generalize
- Three kinds of error
 - Inherent: unavoidable
 - Bias: due to oversimplifications
 - Variance: due to inability to perfectly estimate parameters from limited data



What to remember about classifiers

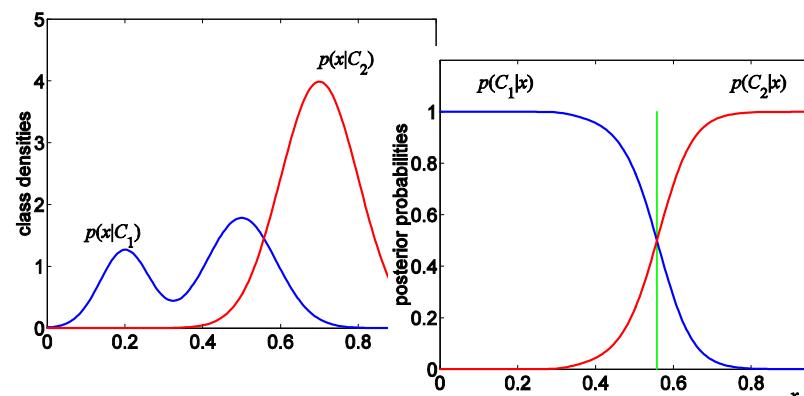
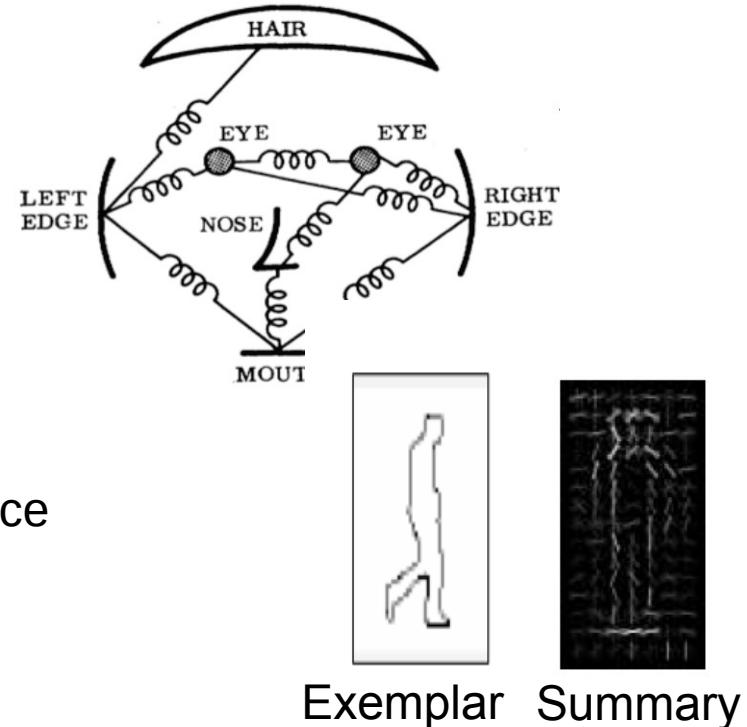
- Machine learning algorithms are tools, not dogmas
- Try simple classifiers first
- Better to have smart features and simple classifiers than simple features and smart classifiers
- Use increasingly powerful classifiers with more training data (bias-variance tradeoff)

How to reduce variance?

- Choose a simpler classifier
- Regularize the parameters
- Get more training data

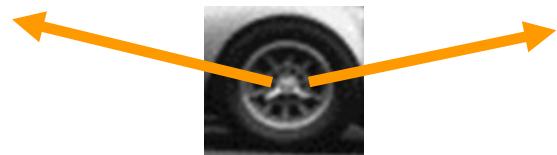
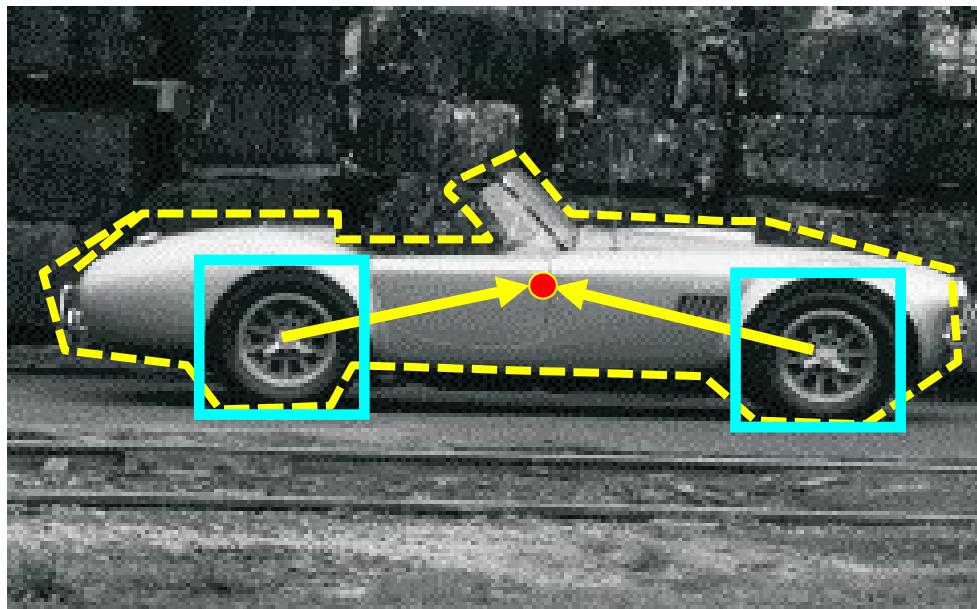
Review: Typical Components

- **Hypothesis** generation
 - Sliding window, Segmentation, feature point detection, random, search
- **Encoding** of (local) image data
 - Colors, Edges, Corners, Histogram of Oriented Gradients, Wavelets, Convolution Filters
- **Relationships** of different parts to each other
 - Blur or histogram, Tree/Star, Pairwise/Covariance
- **Learning** from labeled examples
 - Selecting representative examples (templates), Clustering, Building a cascade
 - Classifiers: Bayes, Logistic regression, SVM, Decision Trees, AdaBoost, ...
 - Generative vs. Discriminative
- **Verification** - removing redundant, overlapping, incompatible examples
 - Non-Max Suppression, context priors, geometry



Implicit shape models

- Visual codebook is used to index votes for object position



visual codeword with
displacement vectors

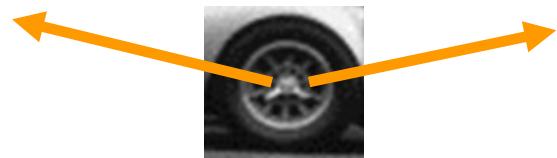
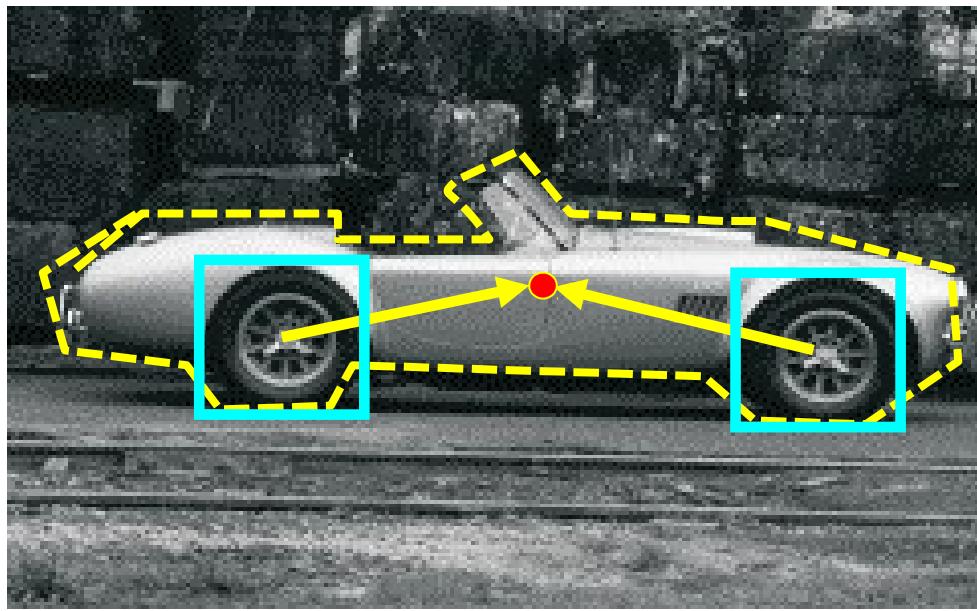
training image annotated with object localization info

B. Leibe, A. Leonardis, and B. Schiele,

[Combined Object Categorization and Segmentation with an Implicit Shape Model](#),
ECCV Workshop on Statistical Learning in Computer Vision 2004

Implicit shape models

- Visual codebook is used to index votes for object position



visual codeword with
displacement vectors

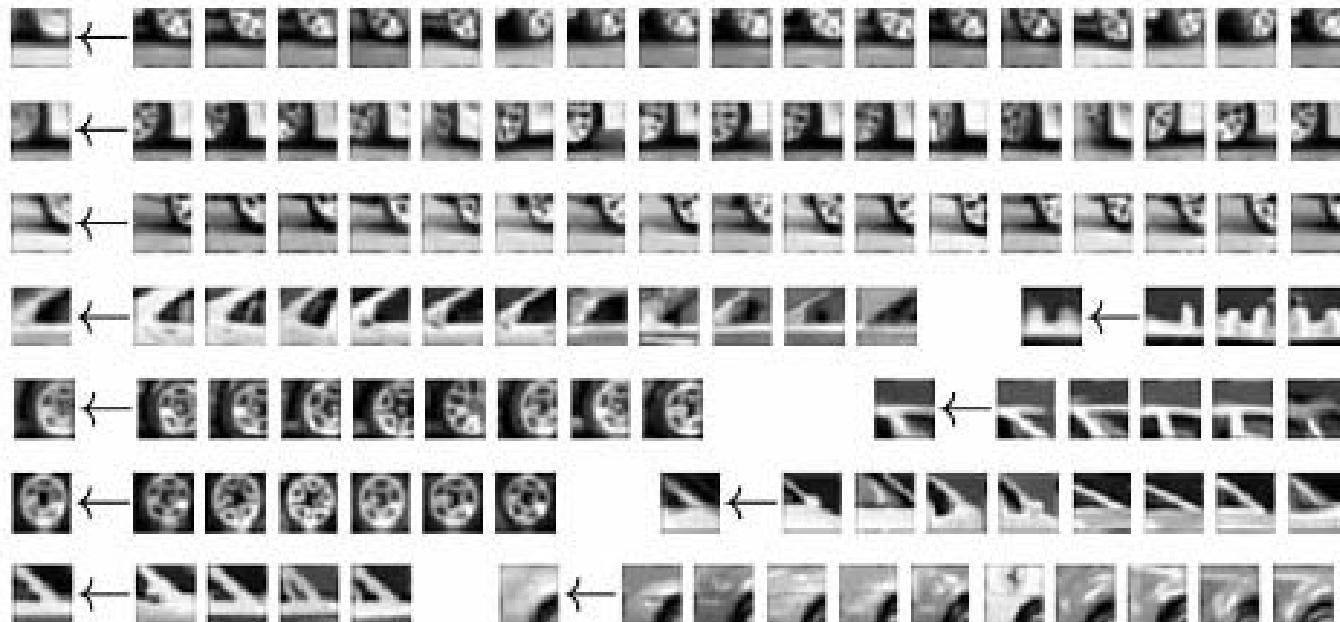
training image annotated with object localization info

B. Leibe, A. Leonardis, and B. Schiele,

[Combined Object Categorization and Segmentation with an Implicit Shape Model](#),
ECCV Workshop on Statistical Learning in Computer Vision 2004

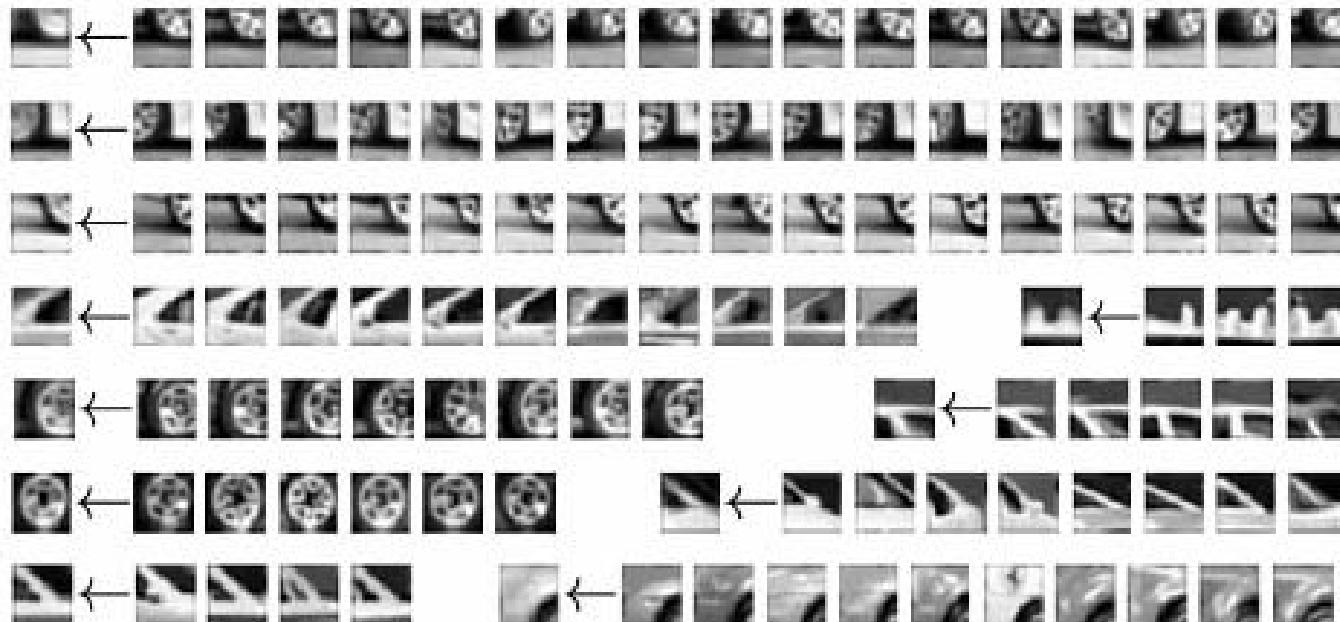
Implicit shape models: Training

1. Build codebook of patches around extracted interest points using clustering



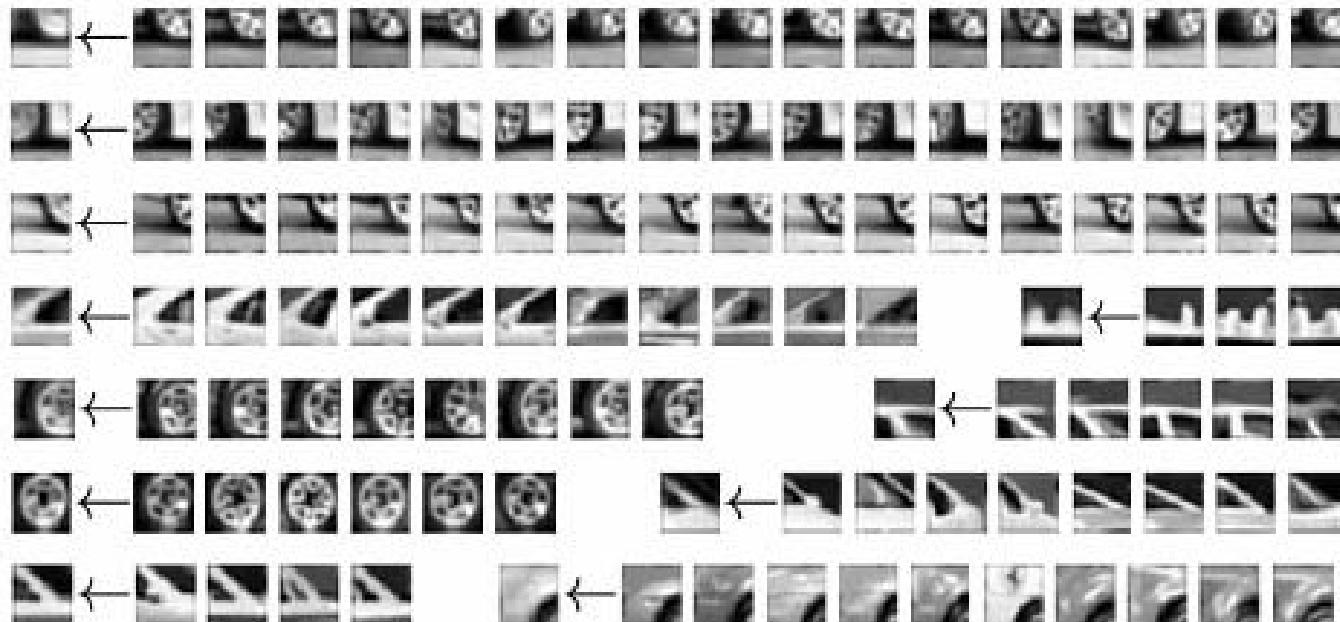
Implicit shape models: Training

1. Build codebook of patches around extracted interest points using clustering



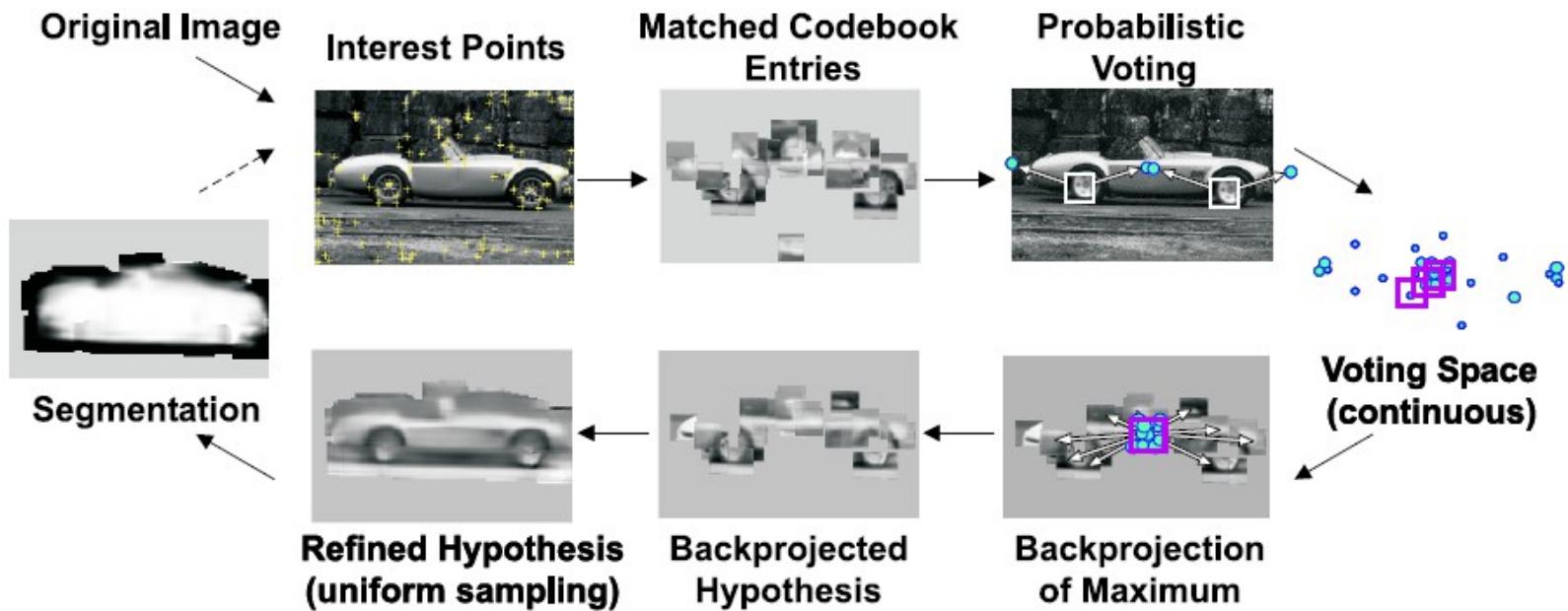
Implicit shape models: Training

1. Build codebook of patches around extracted interest points using clustering



Implicit shape models: Testing

1. Given test image, extract patches, match to codebook entry
2. Cast votes for possible positions of object center
3. Search for maxima in voting space
4. Extract weighted segmentation mask based on stored masks for the codebook occurrences



Example: Results on Cows



Original image

Example: Results on Cows



Interest points

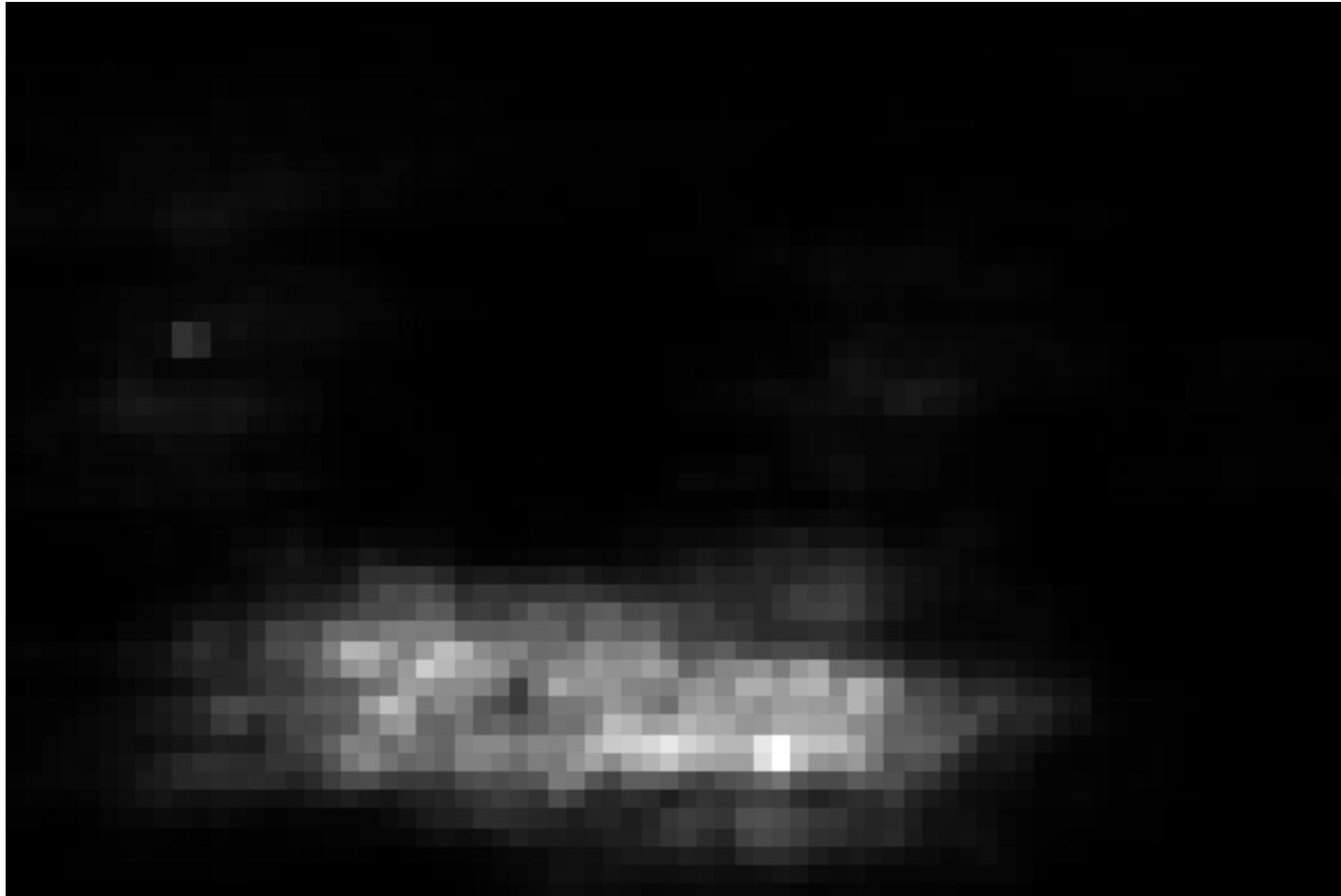
Source: B. Leibe

Example: Results on Cows



Matched patches

Example: Results on Cows



Probabilistic votes

Example: Results on Cows



Hypothesis 1

Example: Results on Cows



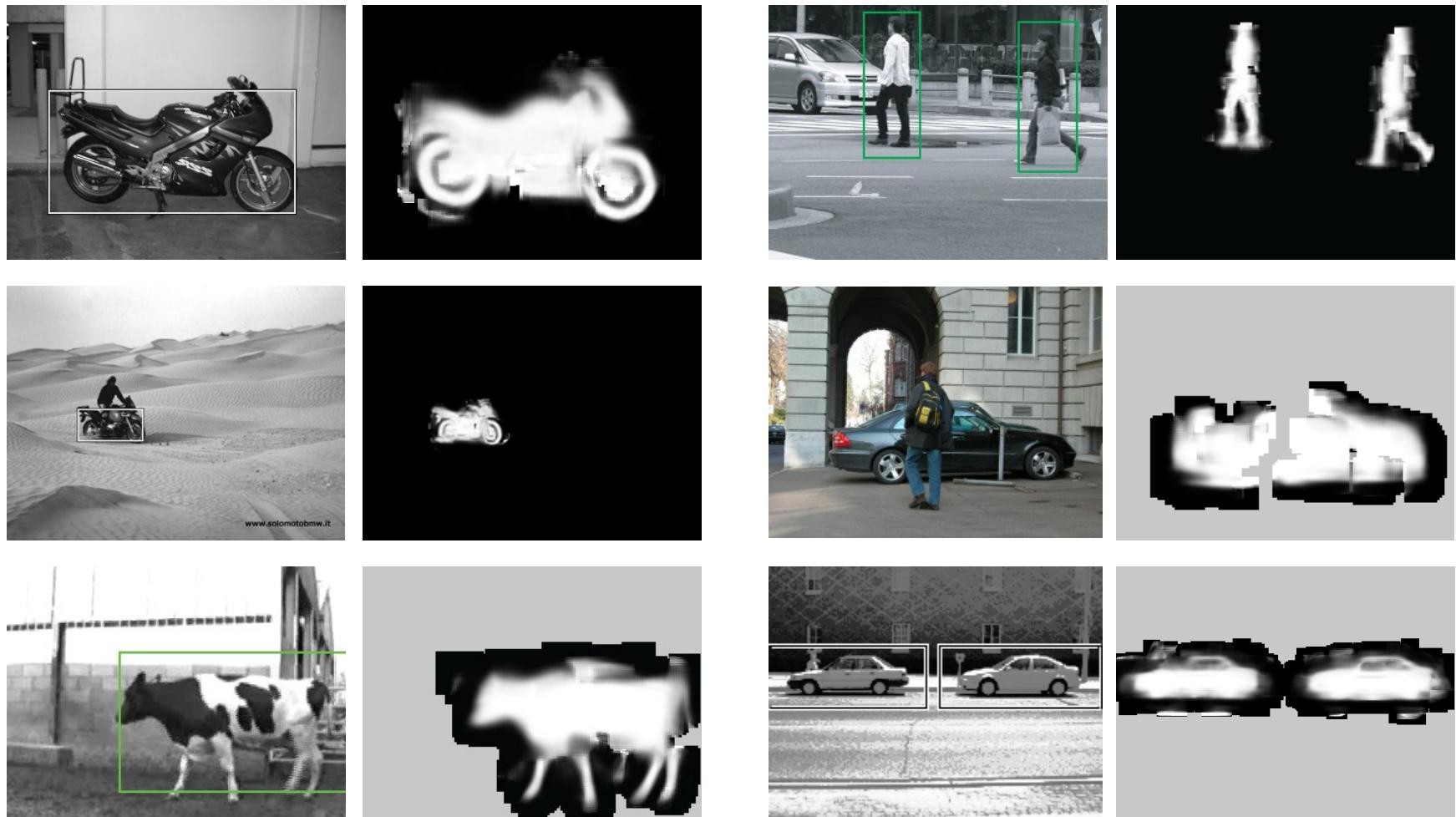
Hypothesis 1

Example: Results on Cows



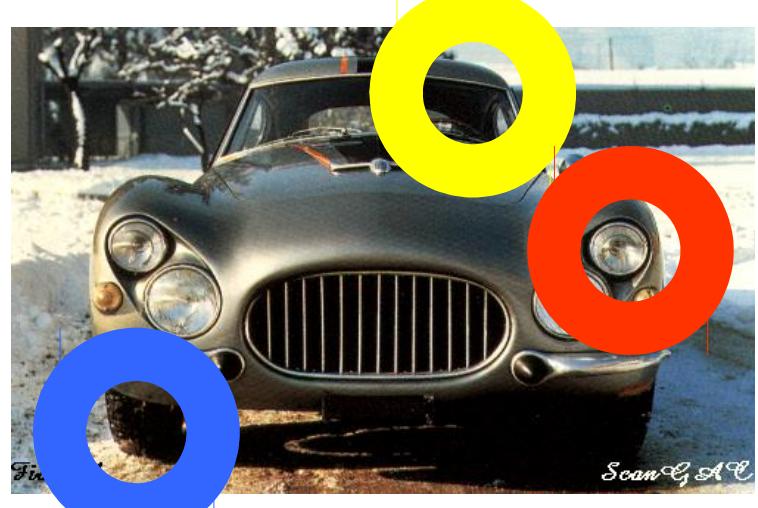
Hypothesis 3

Additional examples



B. Leibe, A. Leonardis, and B. Schiele,
[Robust Object Detection with Interleaved Categorization and Segmentation](#), IJCV
77 (1-3), pp. 259-289, 2008.

Generative part-based models



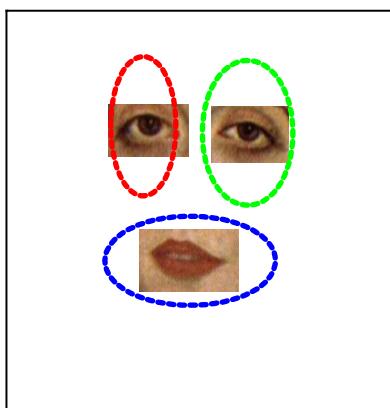
R. Fergus, P. Perona and A. Zisserman,

Object Class Recognition by Unsupervised Scale-Invariant Learning, CVPR 2003

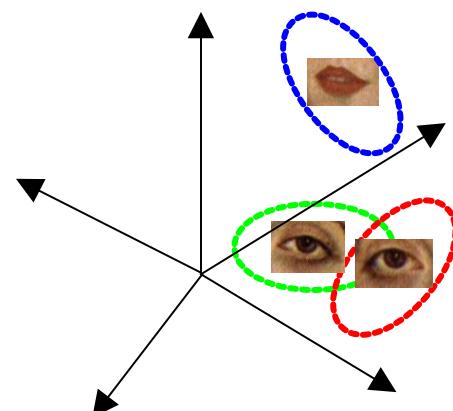
Generative probabilistic model

Foreground model

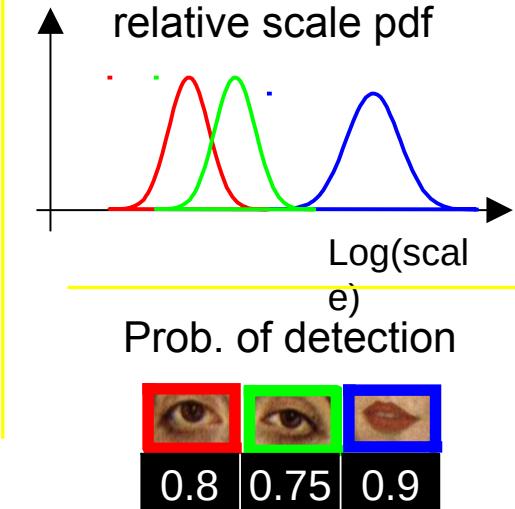
Gaussian shape pdf



Gaussian part appearance pdf

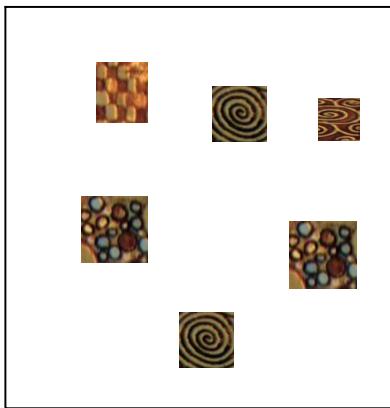


Gaussian relative scale pdf

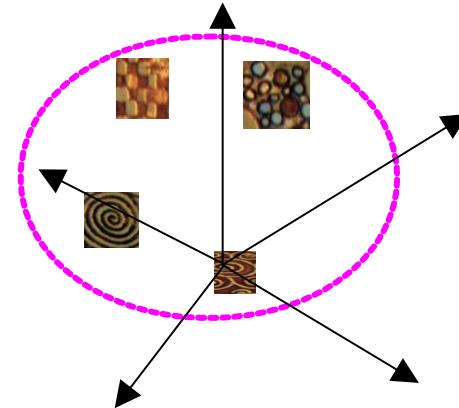


Clutter model

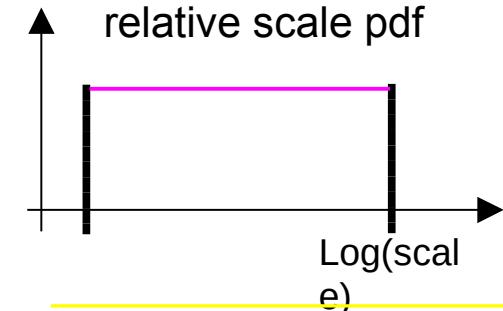
Uniform shape pdf



Gaussian appearance pdf



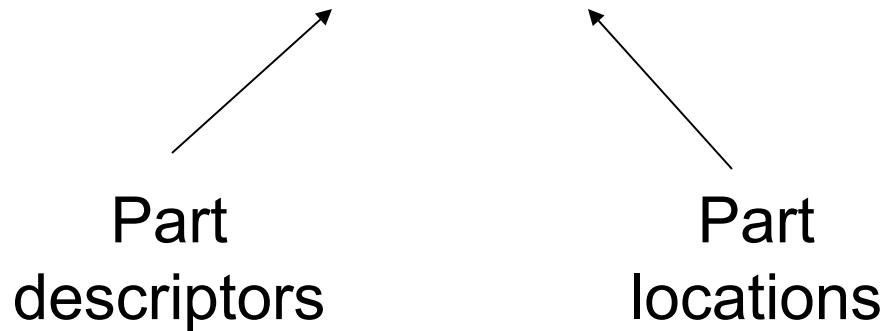
Uniform relative scale pdf



Poisson pdf on # detections

Probabilistic model

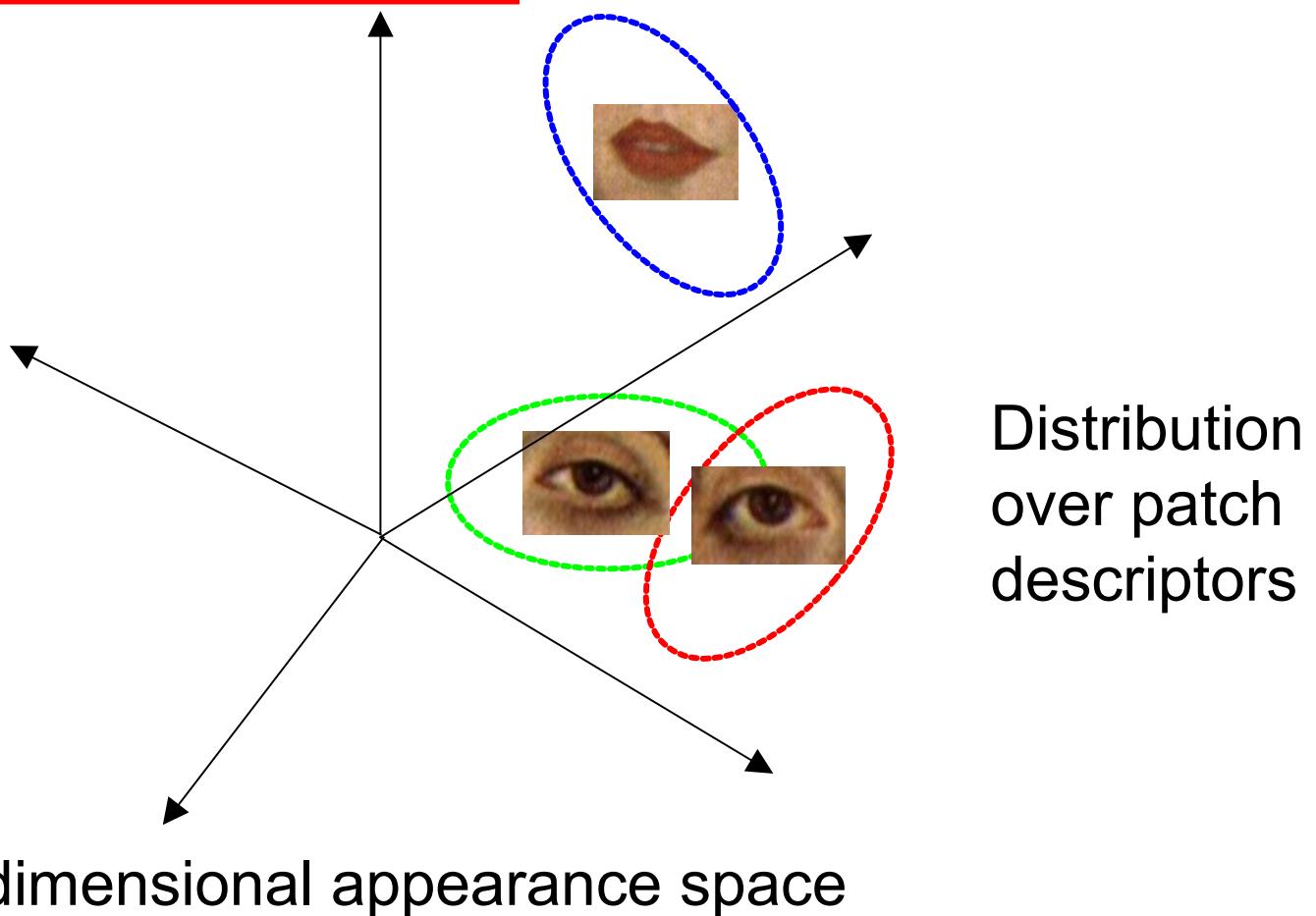
$$P(\text{image} \mid \text{object}) = P(\text{appearance}, \text{shape} \mid \text{object})$$



Candidate parts

Probabilistic model

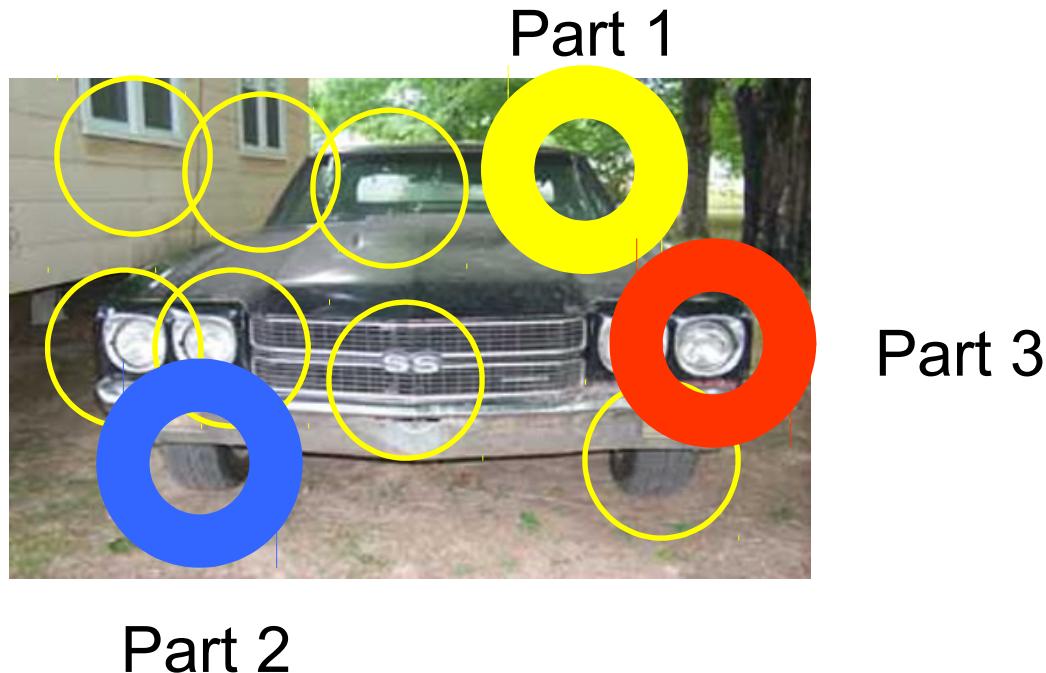
$$\begin{aligned} P(\text{image} \mid \text{object}) &= P(\text{appearance}, \text{shape} \mid \text{object}) \\ &= \max_h P(\text{appearance} \mid h, \text{object}) p(\text{shape} \mid h, \text{object}) p(h \mid \text{object}) \end{aligned}$$



Probabilistic model

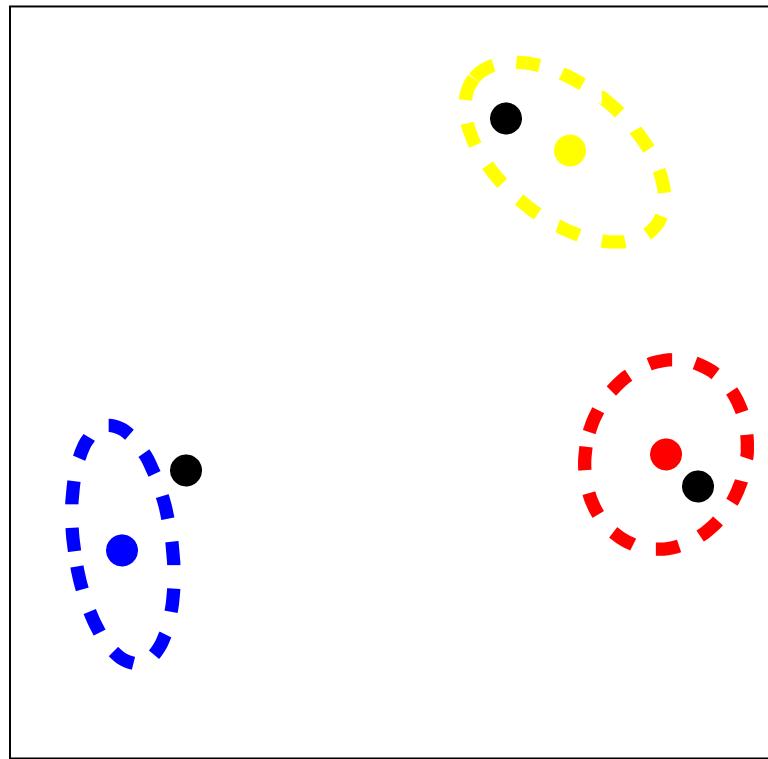
$$\begin{aligned} P(\text{image} \mid \text{object}) &= P(\text{appearance}, \text{shape} \mid \text{object}) \\ &= \max_h P(\text{appearance} \mid h, \text{object}) p(\text{shape} \mid h, \text{object}) p(h \mid \text{object}) \end{aligned}$$

h : assignment of features to parts



Probabilistic model

$$\begin{aligned} P(\text{image} \mid \text{object}) &= P(\text{appearance}, \text{shape} \mid \text{object}) \\ &= \max_h P(\text{appearance} \mid h, \text{object}) p(\text{shape} \mid h, \text{object}) p(h \mid \text{object}) \end{aligned}$$

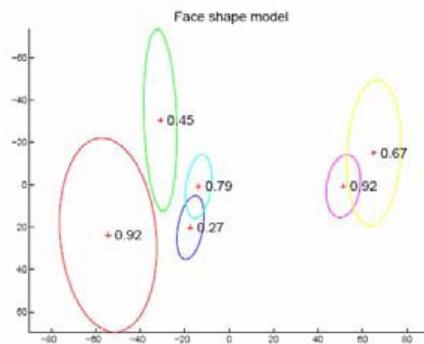


2D image space

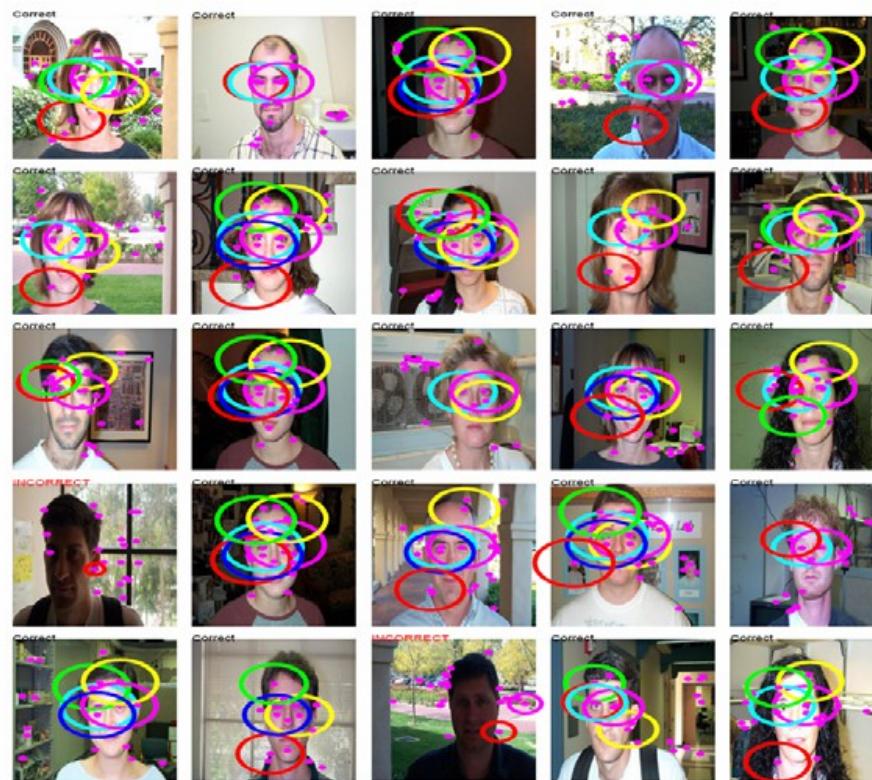
Distribution
over joint
part positions

Results: Faces

Face
shape
model



Recognition
results



Patch
appearance
model

Results: Motorbikes

Part 1 – Det:5e-18



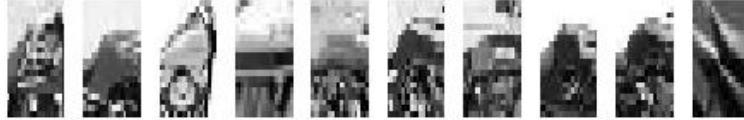
Part 2 – Det:8e-22



Part 3 – Det:6e-18



Part 4 – Det:1e-19



Part 5 – Det:3e-17



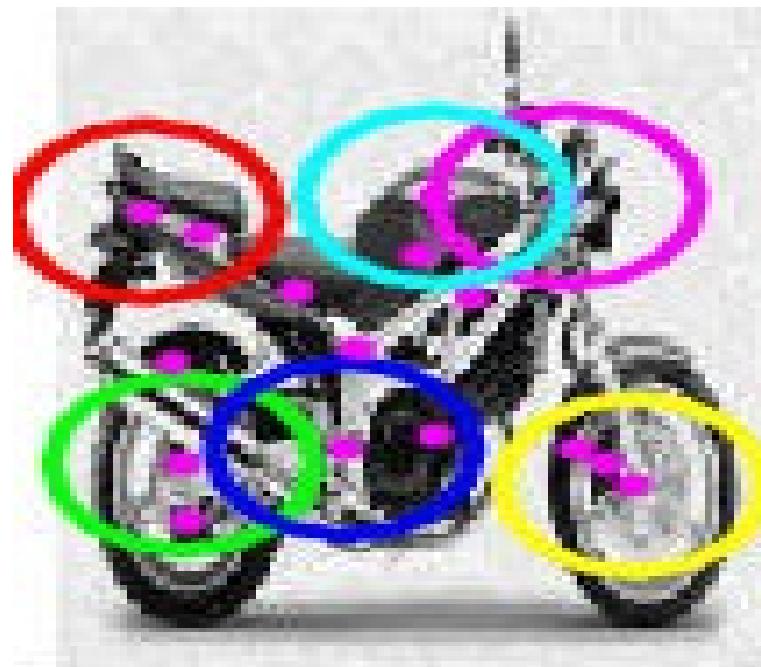
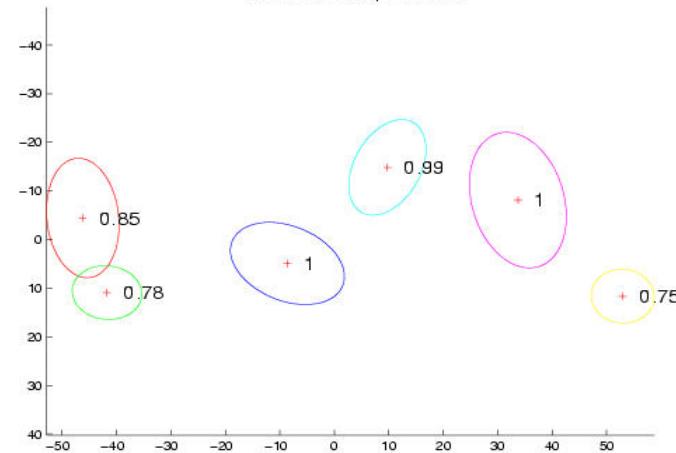
Part 6 – Det:4e-24



Background – Det:5e-19



Motorbike shape model



Motorbikes

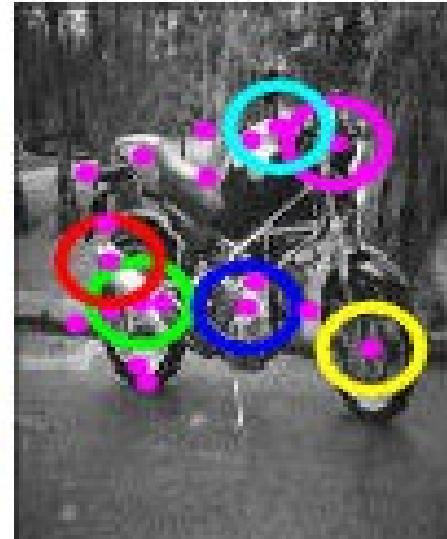
Correct



Correct



Correct



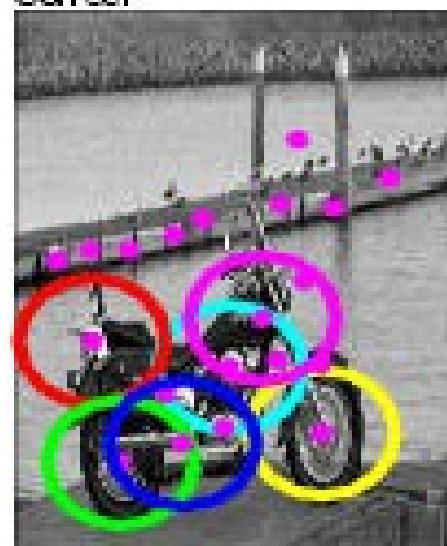
Correct



Correct

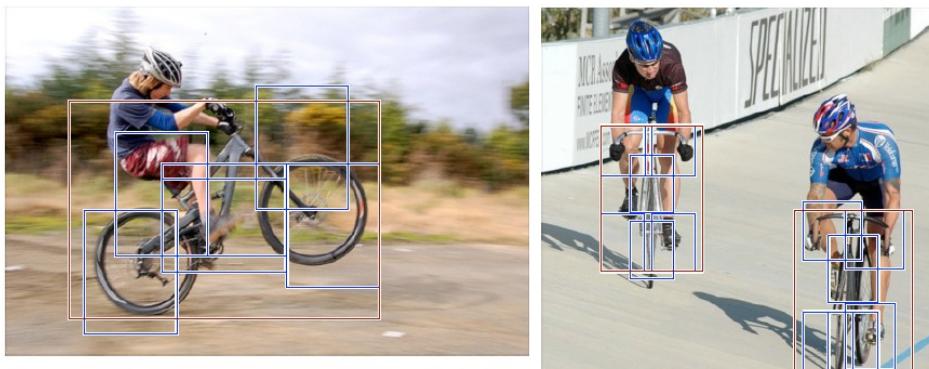


Correct

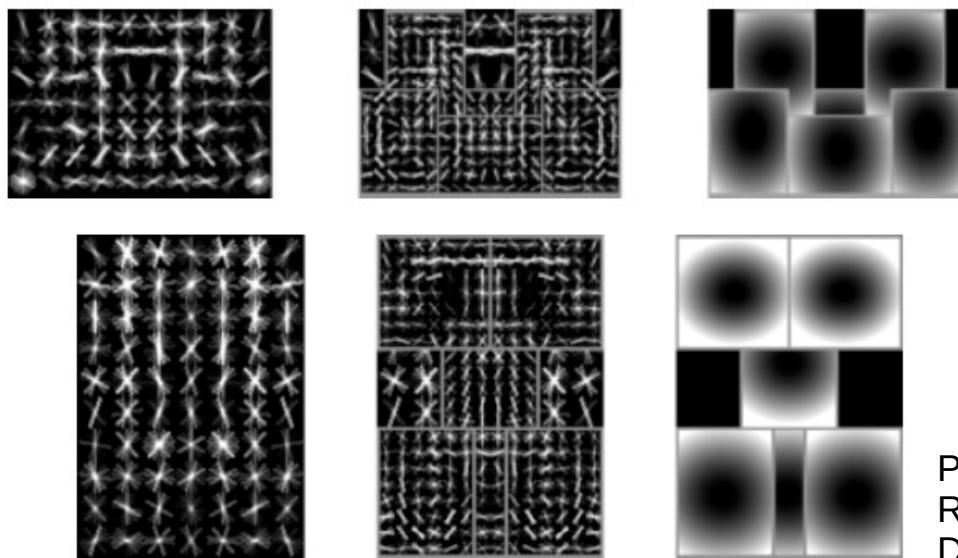


Deformable Parts Model

Detections



Template
Visualization

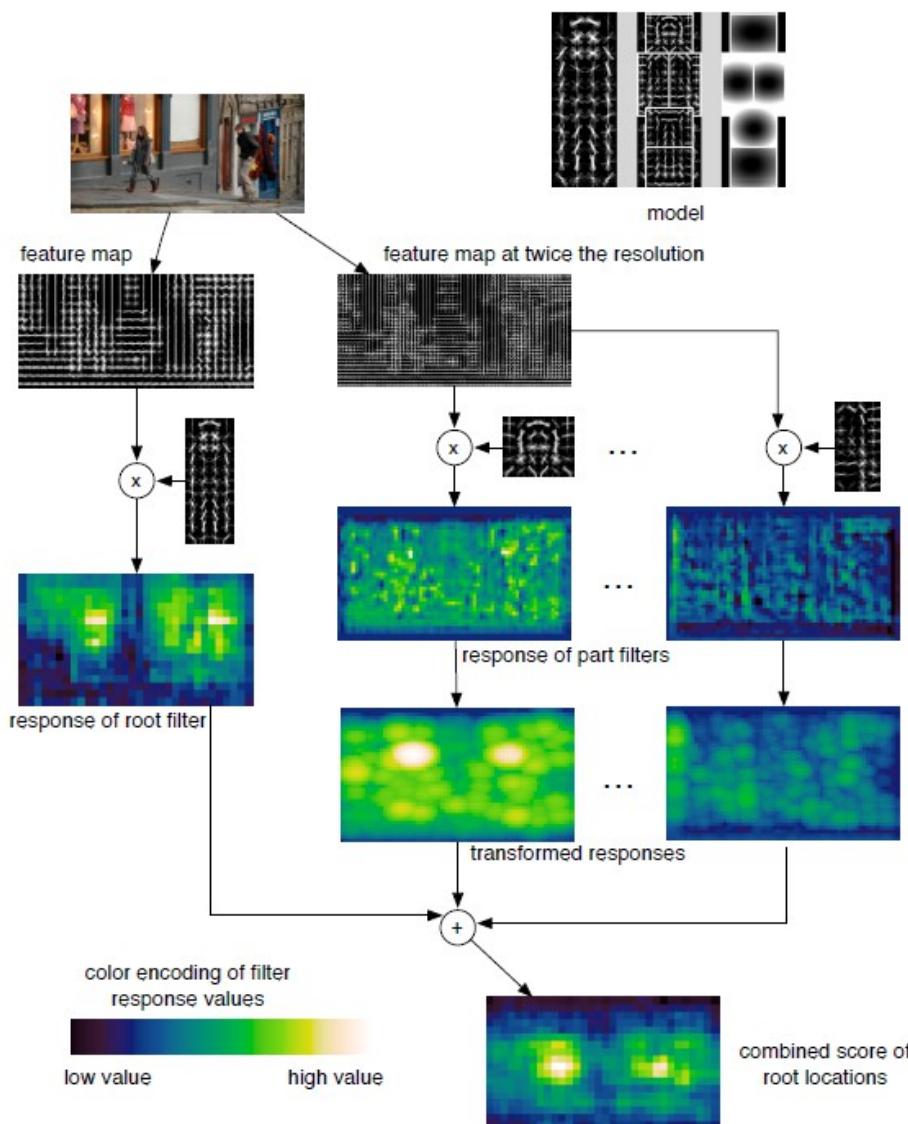


root filters
coarse resolution

part filters
finer resolution

deformation
models

P. Felzenszwalb,
R. Girshick,
D. McAllester,
D. Ramanan
“Object Detection with Discriminatively
Trained Part Based Models”



So many options... How to choose?

- **Hypothesis** generation
 - Sliding window, Segmentation, feature point detection, random, search
- **Encoding** of (local) image data
 - Colors, Edges, Corners, Histogram of Oriented Gradients, Wavelets, Convolution Filters
- **Relationship** of different parts to each other
 - Blur or histogram, Tree/Star, Pairwise/Covariance
- **Learning** from labeled examples
 - Selecting representative examples (templates), Clustering, Building a cascade
 - Classifiers: Bayes, Logistic regression, SVM, Decision Trees, AdaBoost, ...
 - Generative vs. Discriminative
- **Verification** - removing redundant, overlapping, incompatible examples
 - Non-Max Suppression, context priors, geometry

