

מבוא לראייה ממוחשבת – 22928

2016א

מנחה: אמיר אגוזי
egozi5@gmail.com

מפגש מס' 3

מה ראינו בפעם שעברה?

- Model fitting
 - Least squares /total least squares
- Voting
 - Hough transform/generalized Hough transform
- RANSAC

מה היום?

Data analysis •

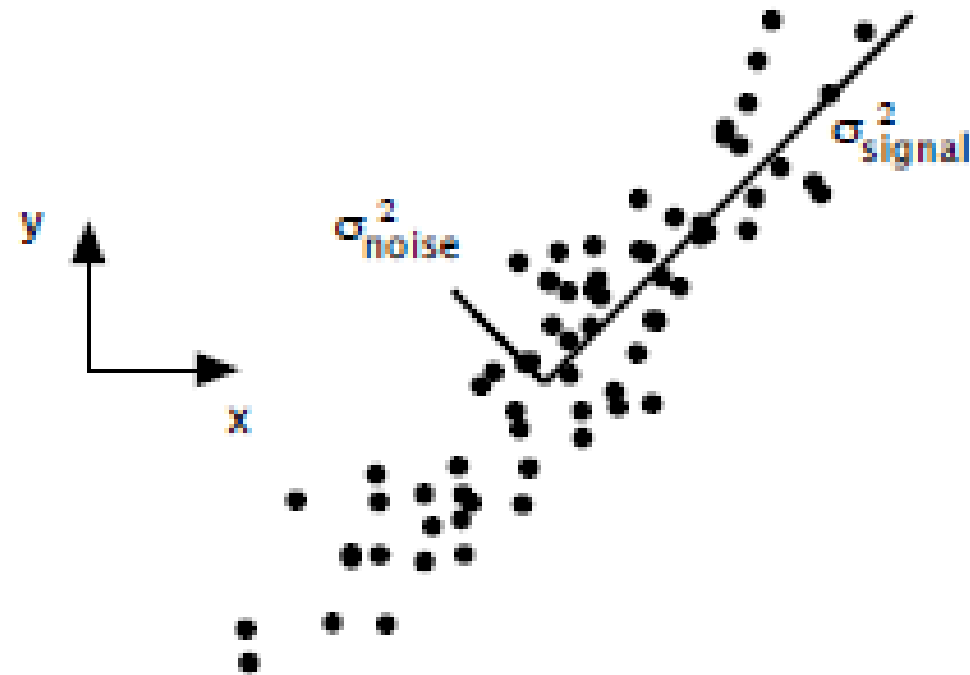
– הורדת מימד לינארית - PCA

– סיווג - Classification

Dimensionality reduction

- Reducing the number of random variables under consideration (**curse of dimensionality**).
- Projection of high-dimensional data to a low-dimensional space that preserves the “important” characteristics of the data.
- Visualize high-dimensional data in low-dimensional space.

Principle component analysis (PCA)



Principle component analysis (PCA)

- Given $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\mathbf{x}_i \in \mathbb{R}^d$
 - Compute sample mean: $\hat{\mu} = \frac{1}{n} \sum_i x_i$
 - Compute sample covariance.

$$\hat{\Sigma} = \frac{1}{n-1} \sum_i (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$$

- Compute eigen-decomposition of $\hat{\Sigma}$

$$\hat{\Sigma} = U\Lambda U^T, \quad \Lambda = (\sigma_1^2, \dots, \sigma_n^2), U^T U = I$$

- Create projection matrix P from
 k “top” eigenvector

$$P = \begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_k^T \end{bmatrix}$$

Principle component analysis (2)

- Given new dataset $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}, z \in \mathbb{R}^d$
- Project according to:

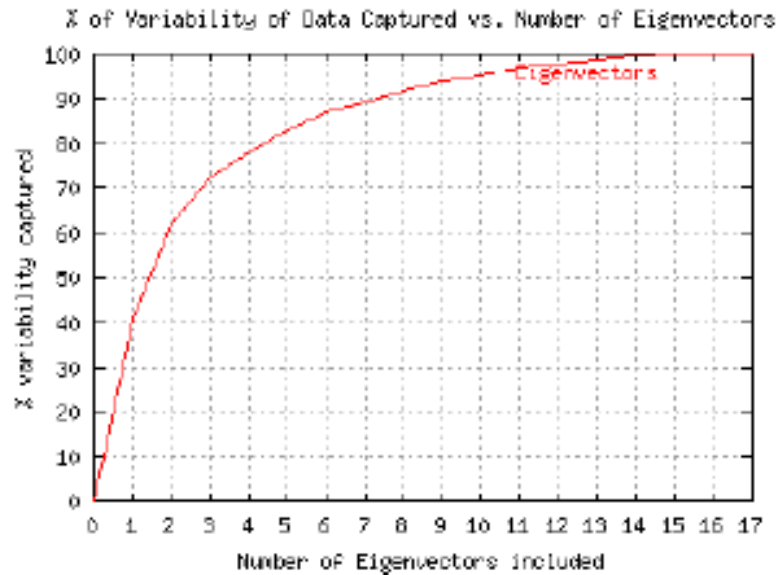
$$\tilde{\mathbf{z}} = P \cdot (\mathbf{z} - \hat{\mu})^T$$

- Use $\tilde{\mathcal{Z}} = \{\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_m\}, z \in \mathbb{R}^k$ to all further processings.

Note the new dimension,
usually: $k \ll d$

Principal component analysis

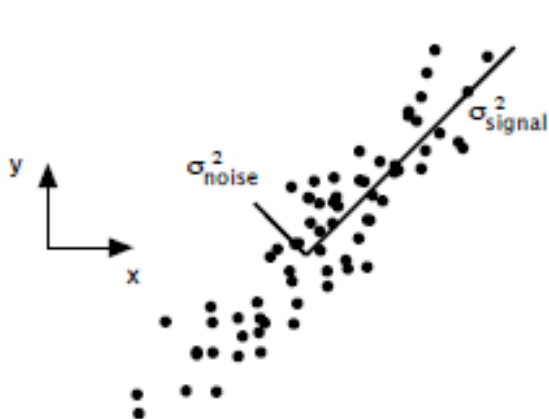
- a natural measure is to pick the eigenvectors that explain X% of the data variability
 - can be done by plotting the ratio r_k as a function of k



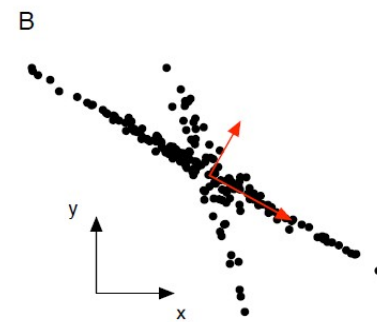
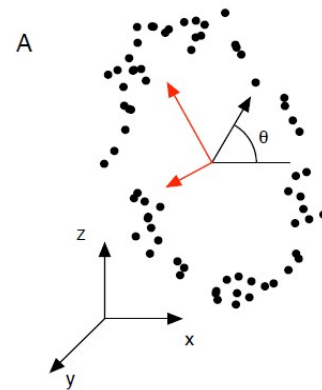
$$r_k = \frac{\sum_{i=1}^k \lambda_i^2}{\sum_{i=1}^n \lambda_i^2}$$

PCA - conclusions

- Simple and intuitive dimensionality reduction approach
- Not always keeping large variance is a good idea.
- Not all datasets exist in a **linear** subspace.
- But, practically...

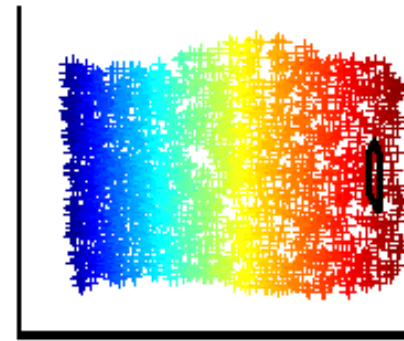
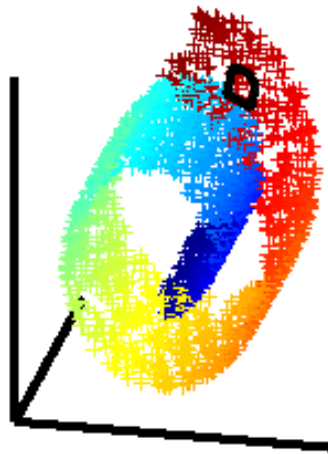
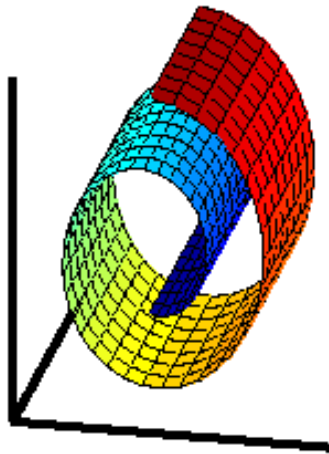


pr
pal



Issues: dimensionality

What if your space isn't *flat*?



Nonlinear methods
LLE, MDS, etc.

Feature extraction

- PCA is a feature extraction technique.
- There are many many more.
- For a review on state-of-the-art techniques:
http://en.wikipedia.org/wiki/Feature_extraction

Classification – Pattern recognition

Subproblems of Pattern classification

- Feature extraction / representation
- Noise
- Over-fitting
- Model selection
- Prior knowledge
- Missing features
- Segmentation
- Context
- Invariances
- Cost and risks
- Computational complexity

Machine learning

- **Supervised learning** - Learn from examples, ((input, output) pairs)
 - **Regression**: predict a continuous value as a function of the input.
 - **Classification**: predict the class label of the input object
- **Unsupervised learning** - fit a model to observations without a priori output.
 - Clustering
 - Identifying patterns in the data.

Unsupervised Learning

- Input – unlabeled data samples $\{\mathbf{x}_i\}, i = 1, \dots, n$
- Why study unlabeled data?
 - Labeled data can be costly
 - Cluster first, label later
 - Identify features that will be useful for categorization
 - Exploratory data analysis

Supervised learning

Recognition: A machine learning approach



Slides adapted from Fei-Fei Li, Rob Fergus, Antonio Torralba, Kristen Grauman, and Derek Hoiem

The machine learning framework

- Apply a prediction function to a feature representation of the image to get the desired output:

$f(\text{apple image}) = \text{"apple"}$

$f(\text{tomato image}) = \text{"tomato"}$

$f(\text{cow image}) = \text{"cow"}$

The machine learning framework

$$y = f(\mathbf{x})$$

A diagram illustrating the machine learning equation $y = f(\mathbf{x})$. The equation is written in blue. Below it, three red arrows point upwards to the components: the first arrow points to y and is labeled 'output'; the second arrow points to f and is labeled 'prediction function'; the third arrow points to \mathbf{x} and is labeled 'Image feature'.

Training: given a *training set* of labeled examples

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

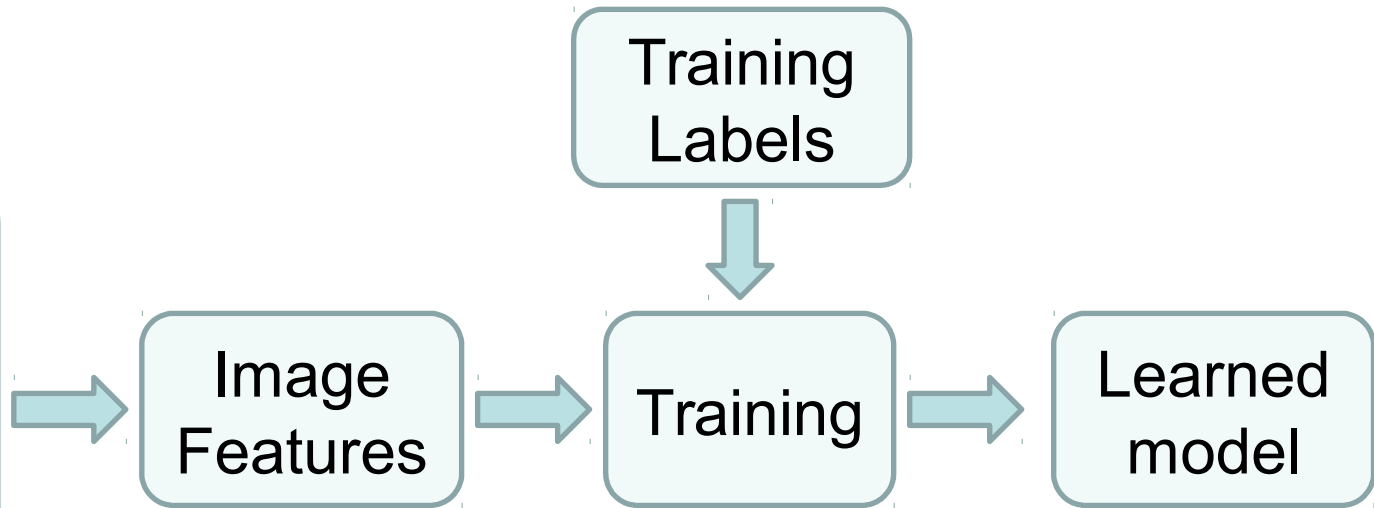
estimate the prediction function $f(\mathbf{x})$ by minimizing the prediction error on the training set

- **Testing:** apply f to a never before seen *test example* \mathbf{x} and output the predicted value $y = f(\mathbf{x})$

Step I - Learning

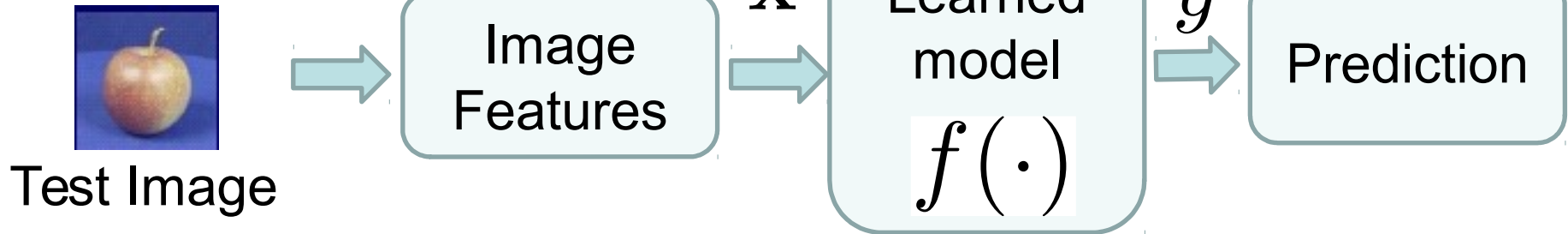
Training

Training
Images



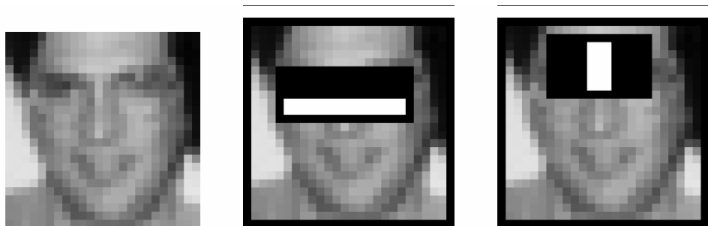
Step II - Test

Testing

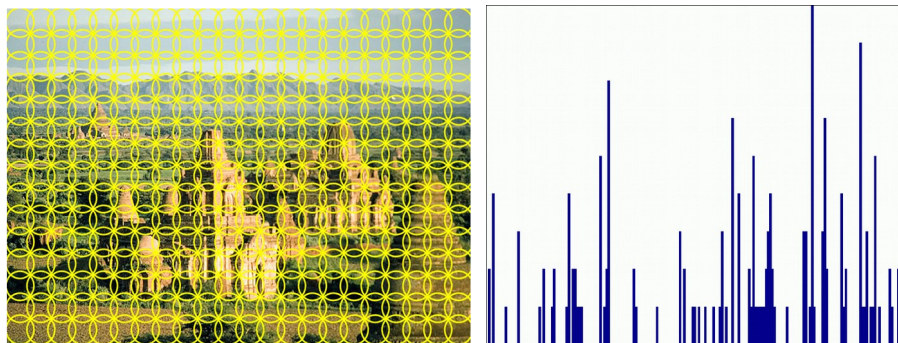


Features (examples)

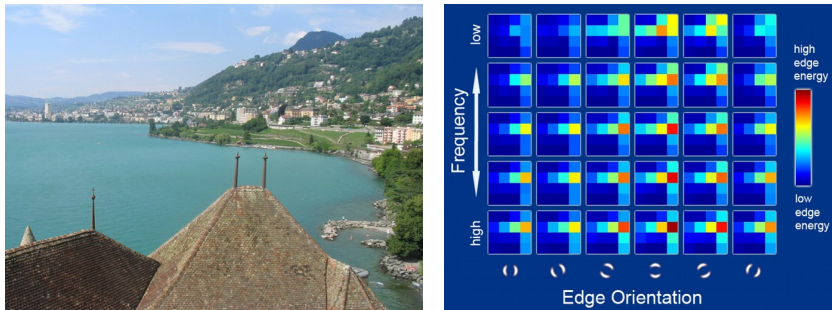
- Raw pixels (and simple functions of raw pixels)



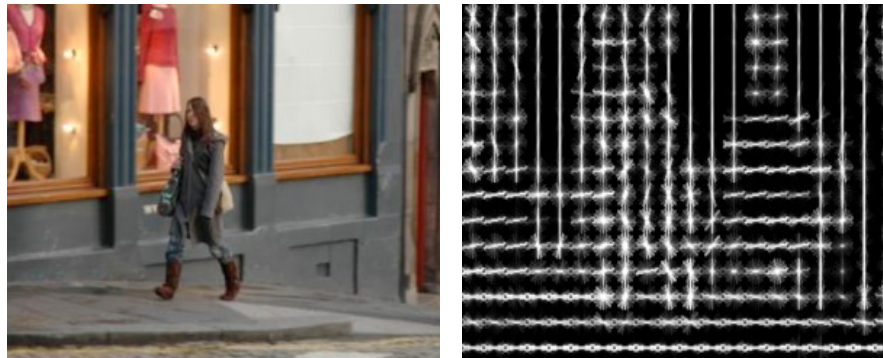
- Histograms, bags of features



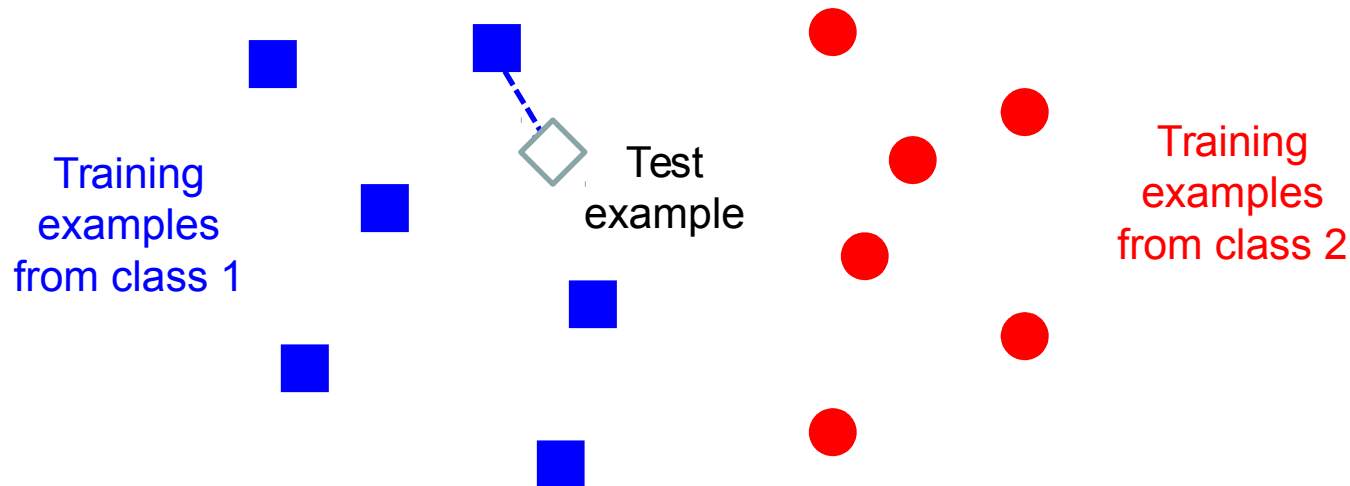
- GIST descriptors



- Histograms of oriented gradients (HOG)



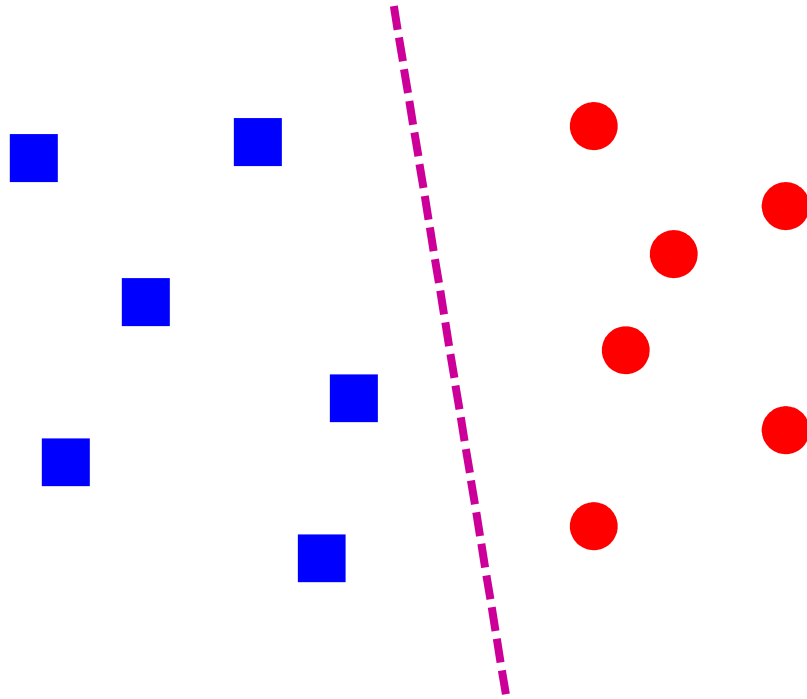
Classifiers: Nearest neighbor



$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

- All we need is a distance function for our inputs
- No training required!

Classifiers: Linear



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

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

Generalization



Training set (labels known)



Test set (labels unknown)

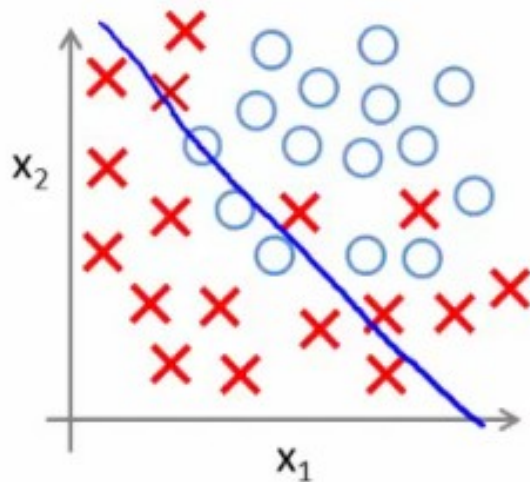
- How well does a learned model *generalize* from the data it was trained on to a new test set?

Diagnosing generalization ability

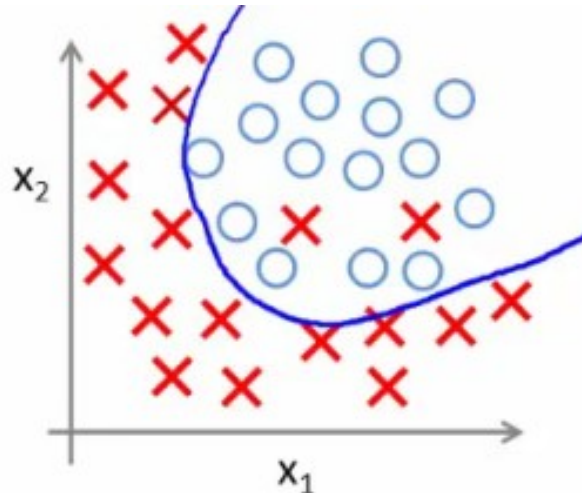
- **Training error:** how well does the model perform at prediction on the data on which it was trained?
- **Test error:** how well does it perform on a never before seen test set?
- Training and test error are both *high*: **underfitting**
 - Model does an equally poor job on the training and the test set
 - Either the training procedure is ineffective or the model is too “simple” to represent the data
- Training error is *low* but test error is *high*: **overfitting**
 - Model has fit irrelevant characteristics (noise) in the training data
 - Model is too complex or amount of training data is insufficient

Underfitting and overfitting

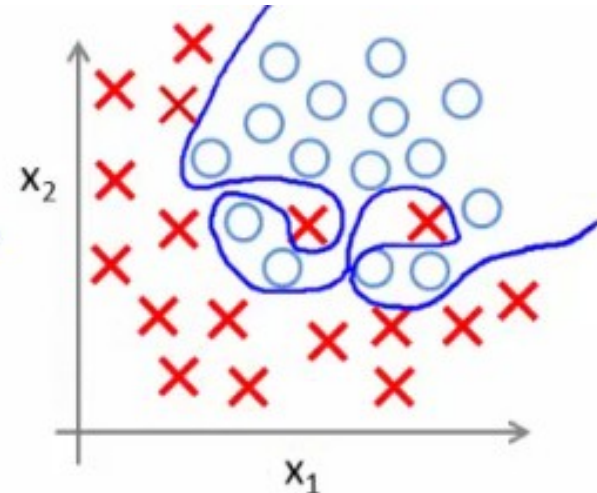
Underfitting



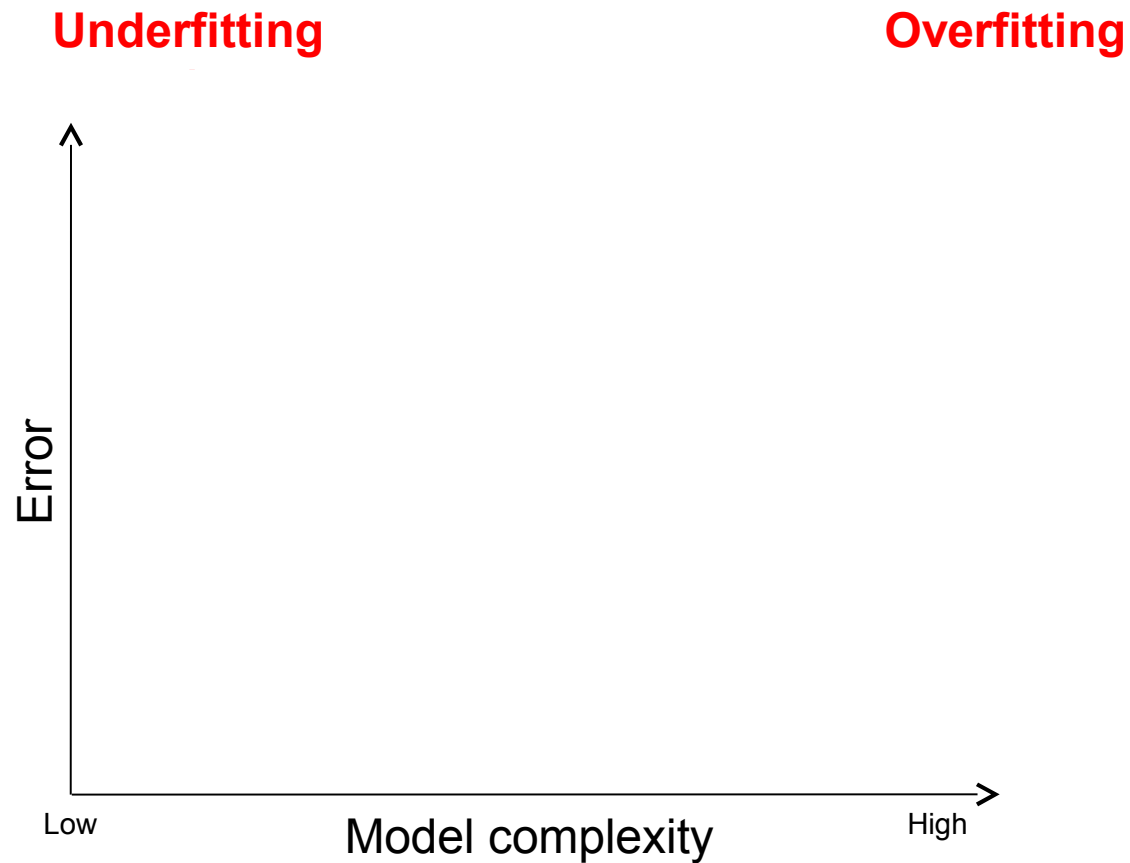
Good generalization



Overfitting

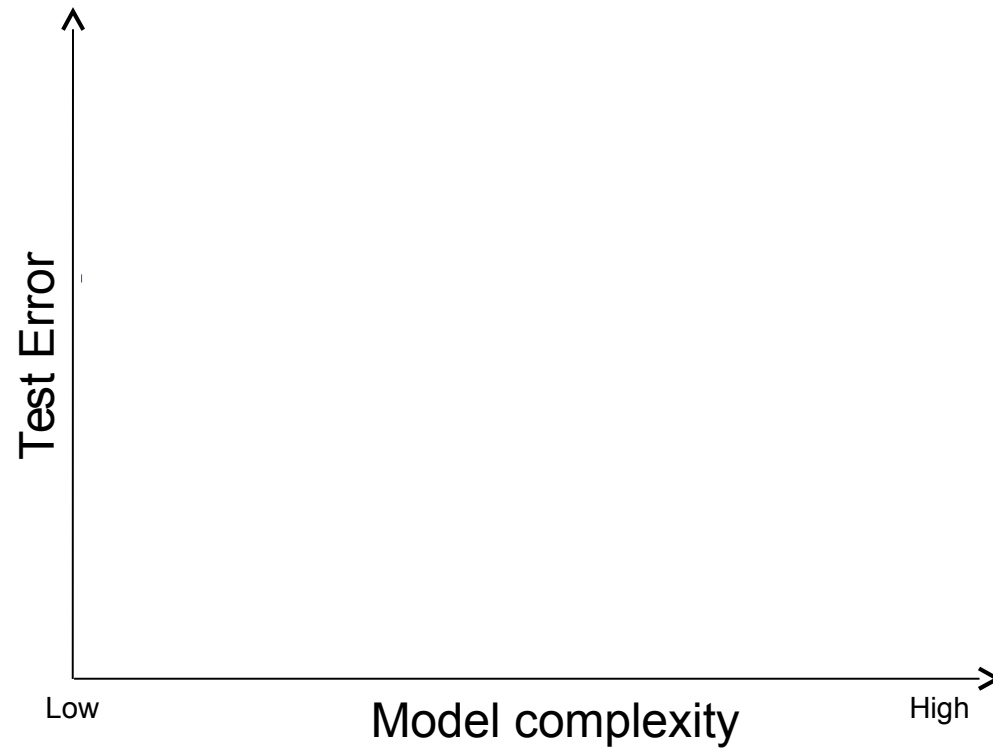


Effect of model complexity



Slide credit: D. Hoiem

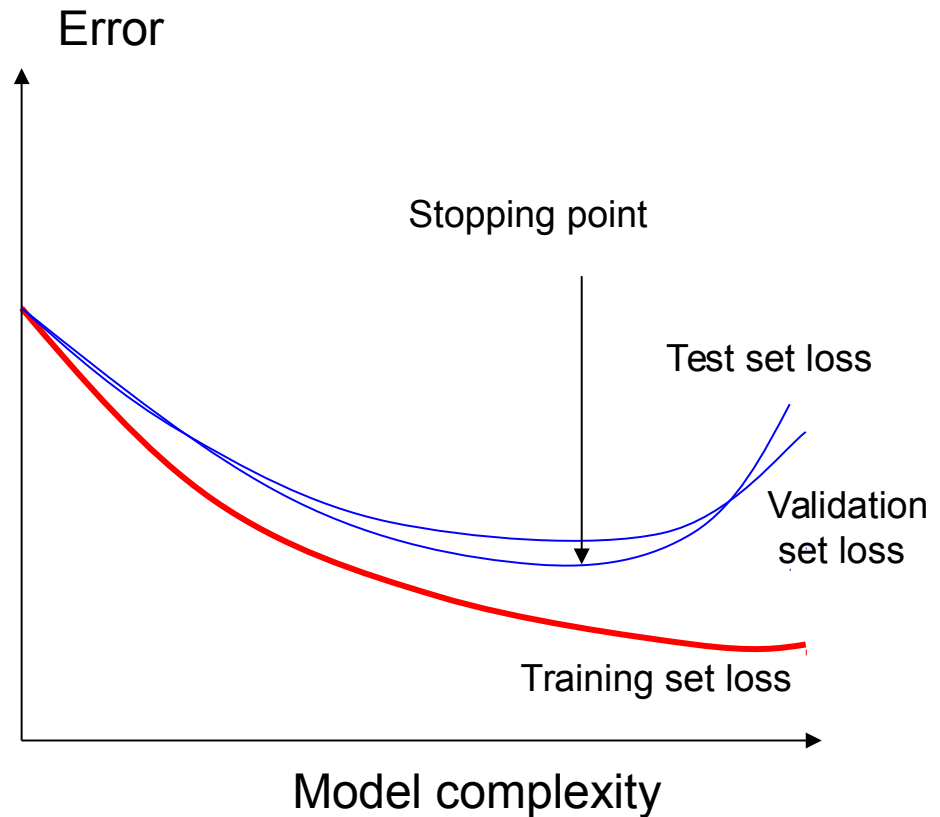
Effect of training set size



Slide credit: D. Hoiem

Validation

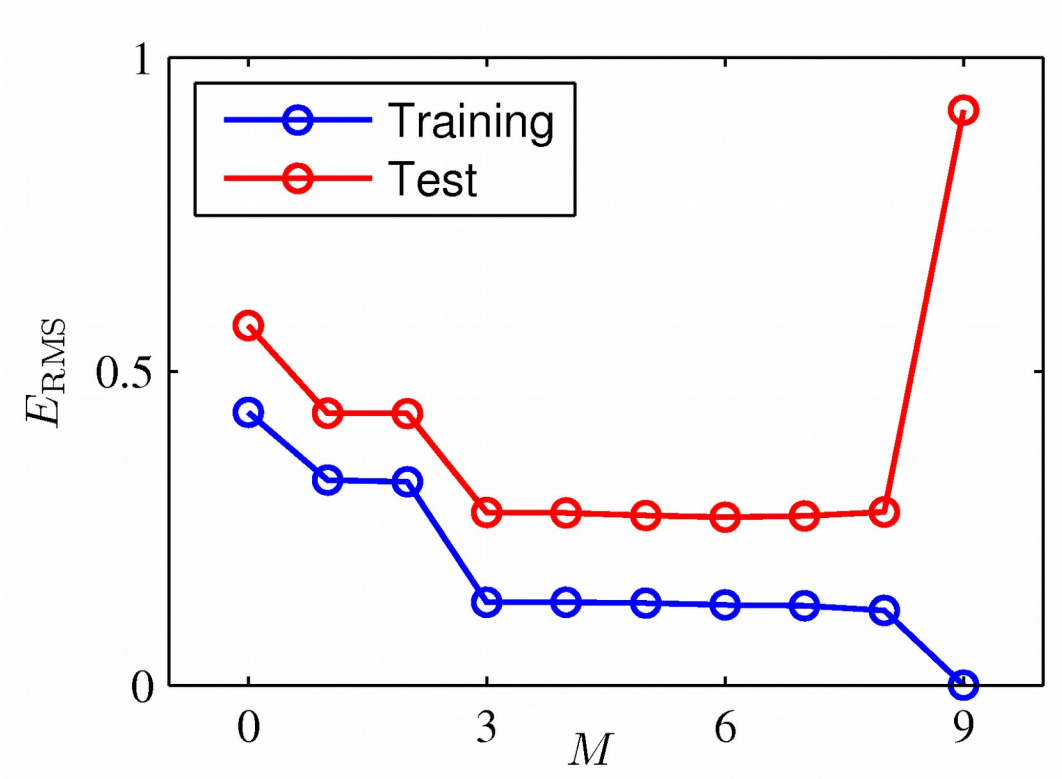
- Split the dataset into **training, validation, and test** sets
- Use training set to **optimize model parameters**
- Use validation test to **choose the best model**
- Use test set only to **evaluate performance**



Cross validation

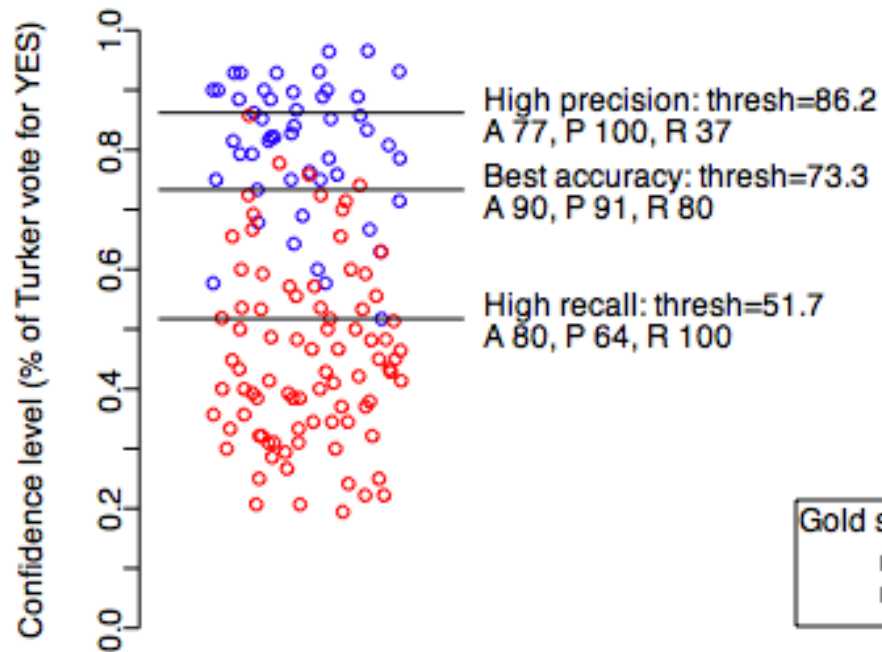
	Downside	Upside
Test-set	Variance: unreliable estimate of future performance	Cheap
Leave-one-out	Expensive. Has some weird behavior	Doesn't waste data
10-fold	Wastes 10% of the data. 10 times more expensive than test set	Only wastes 10%. Only 10 times more expensive instead of R times.
3-fold	Wastier than 10-fold. Expensivier than test set	Slightly better than test-set
R-fold	Identical to Leave-one-out	

Generalization



Classifier performance

Test set separation by Turker ensemble binary classifier



**prediction
outcome**

Positive

Negative

actual value

Positive

Negative

True
Positive

False
Positive

False
Negative

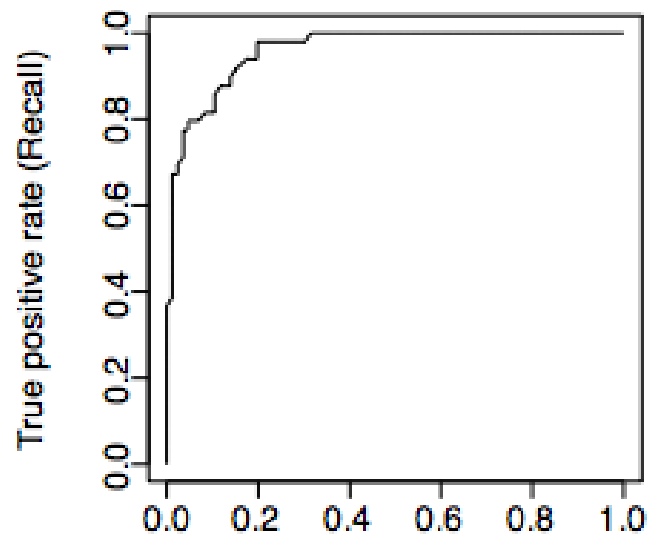
True
Negative

Above a threshold, classified as Y, below as N
Errors above are false pos; errors below are false neg
Accuracy, Precision, Recall in %
Dots have horizontal jitter (x-axis has no meaning)

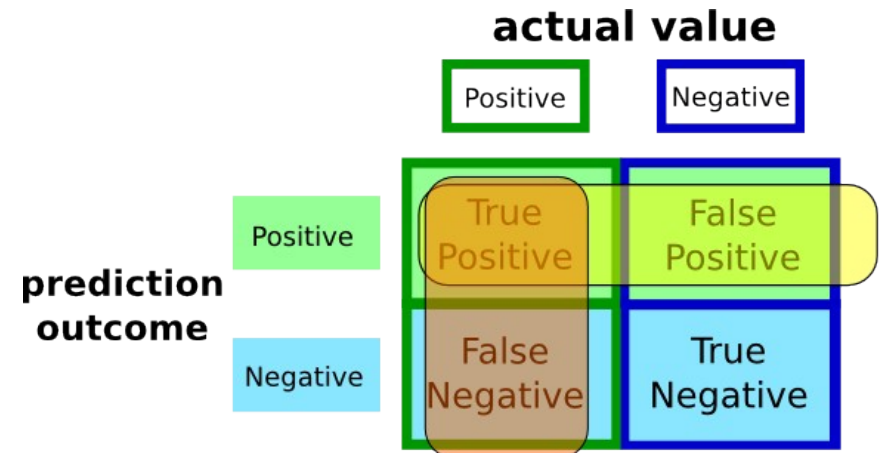
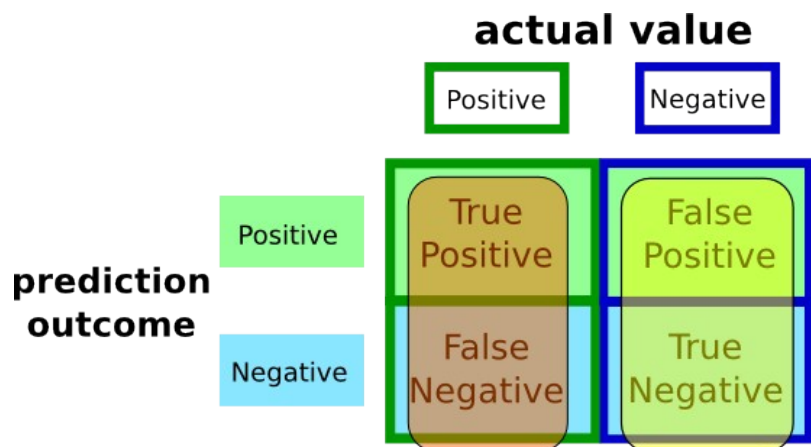
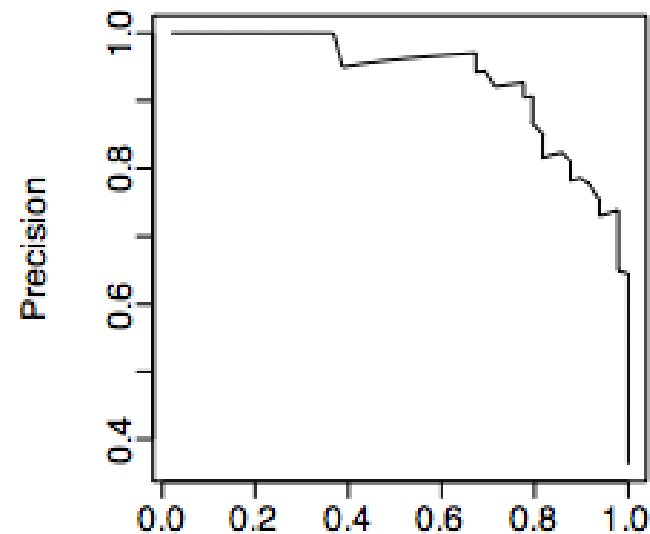
<http://blog.doloreslabs.com/?p=61>

Classifier performance

ROC curve



Precision-Recall curve



Object categorization: the statistical viewpoint

- MAP decision: $p(\text{zebra} \mid \text{image})$
vs.
 $p(\text{no zebra} \mid \text{image})$



Object categorization: the statistical viewpoint

- MAP decision: $p(\text{zebra} \mid \text{image})$
vs.
 $p(\text{no zebra} \mid \text{image})$



- Bayes rule:

$$\underbrace{p(\text{zebra} \mid \text{image})}_{\text{posterior}} \propto \underbrace{p(\text{image} \mid \text{zebra})}_{\text{likelihood}} \underbrace{p(\text{zebra})}_{\text{prior}}$$

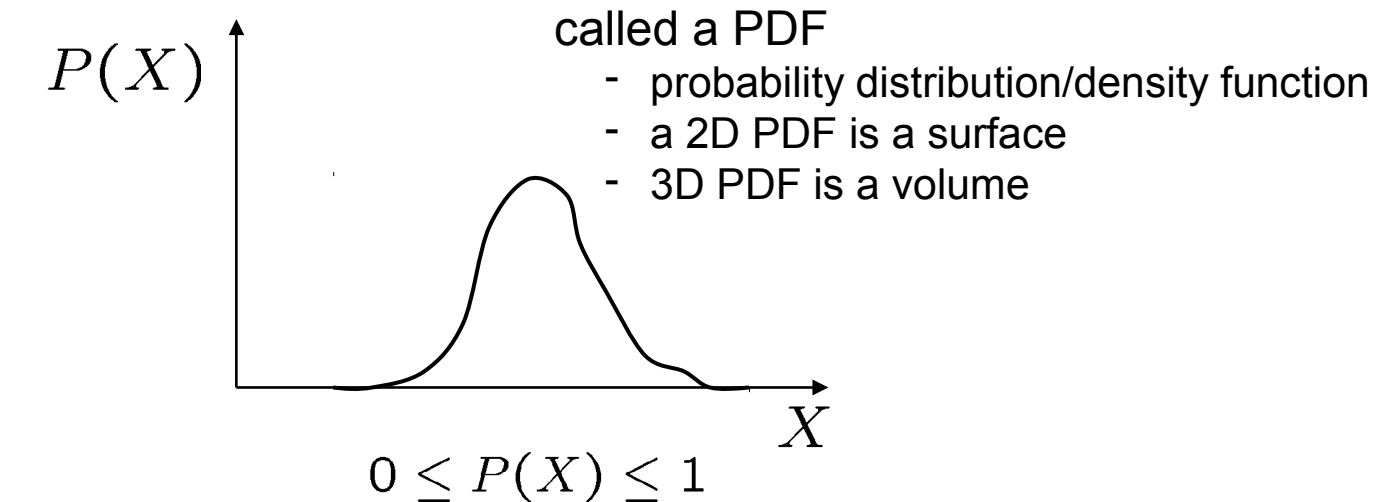
Object categorization: the statistical viewpoint

$$\underbrace{p(\textit{zebra} | \textit{image})}_{\text{posterior}} \propto \underbrace{p(\textit{image} | \textit{zebra})}_{\text{likelihood}} \underbrace{p(\textit{zebra})}_{\text{prior}}$$

- **Discriminative methods:** model posterior
- **Generative methods:** model likelihood and prior

Probability

- X is a random variable.
- $P(X)$ is the probability that X achieves a certain value.



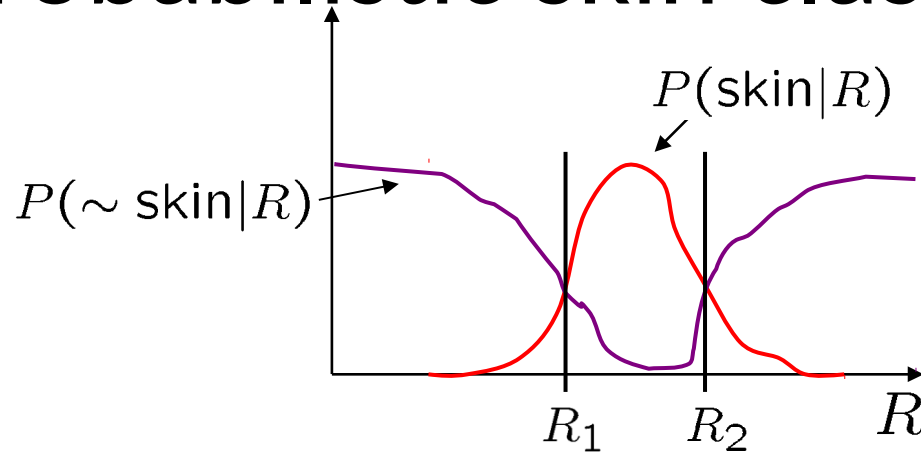
$$\int_{-\infty}^{\infty} P(X) dX = 1$$

continuous X

$$\sum P(X) = 1$$

discrete X

Probabilistic skin classification



- Model PDF / uncertainty
 - Each pixel has a probability of being skin or not skin
$$P(\sim \text{skin}|R) = 1 - P(\text{skin}|R)$$
- Skin classifier
 - Given $X = (R,G,B)$: how to determine if it is skin or not?
 - Choose interpretation of highest probability
- Where do we get $P(\text{skin}|R)$ $P(\sim \text{skin}|R)$

התפלגות גאוסיאנית רב-ממדית

- בהנתן:

Σ – (מטריצת שונות משותפת)

μ – (ממוצע רב-ממדי)

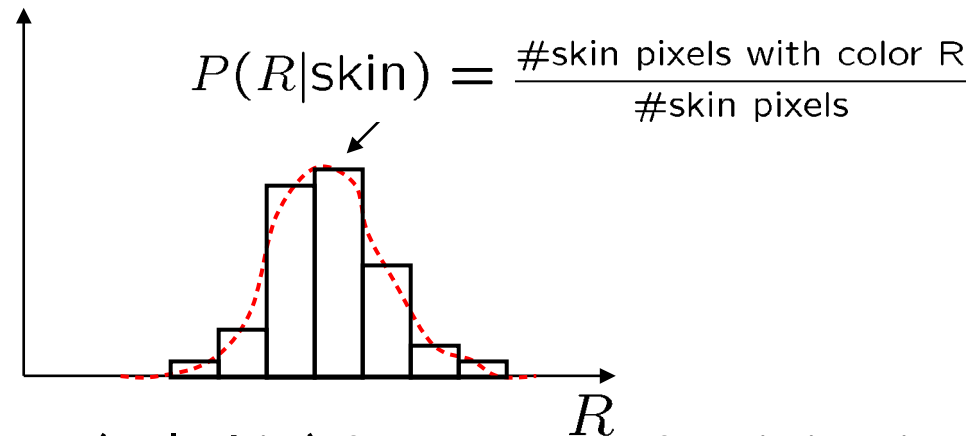
- k מממד x חישוב הסתברות לנקודה

$$P(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

- MATLAB - או ב

`y = mvnpdf(X,MU,SIGMA);`

Learning conditional PDF's



- We can calculate $P(R | \text{skin})$ from a set of training images
- But this isn't quite what we want
 - We want $P(\text{skin} | R)$ not $P(R | \text{skin})$
 - How can we get it?

Bayes rule

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

what we measure
(likelihood)

domain knowledge
(prior)

- In terms of our problem:

$$P(\text{skin}|R) = \frac{P(R|\text{skin}) P(\text{skin})}{P(R)}$$

what we want
(posterior)

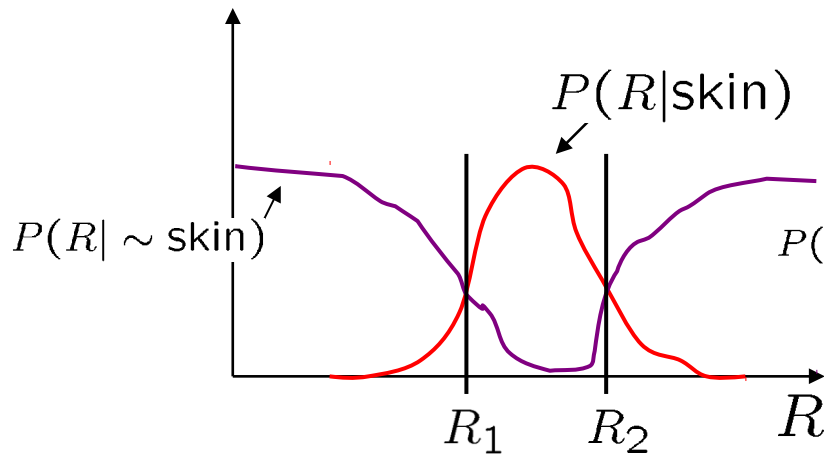
normalization term

$$P(R) = P(R|\text{skin})P(\text{skin}) + P(R|\sim \text{skin})P(\sim \text{skin})$$

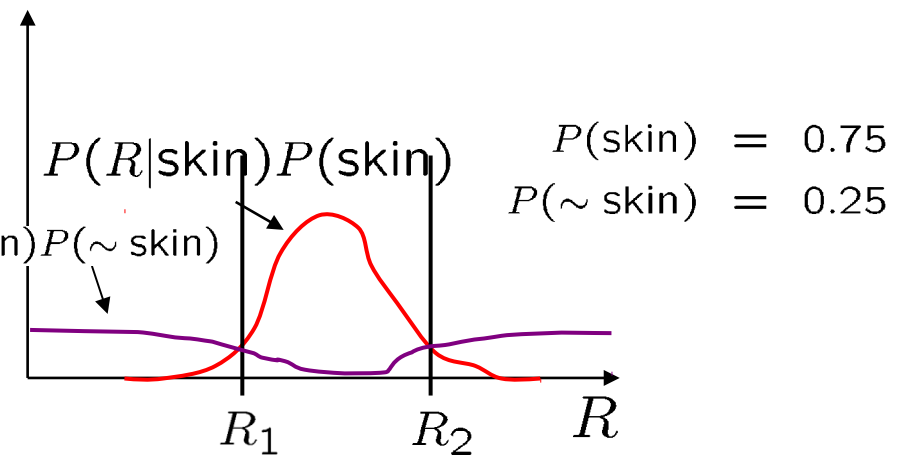
What can we use for the prior $P(\text{skin})$?

- Domain knowledge:
 - $P(\text{skin})$ may be larger if we know the image contains a person
 - For a portrait, $P(\text{skin})$ may be higher for pixels in the center
- Learn the prior from the training set. How?
 - $P(\text{skin})$ is proportion of skin pixels in training set

Bayesian estimation



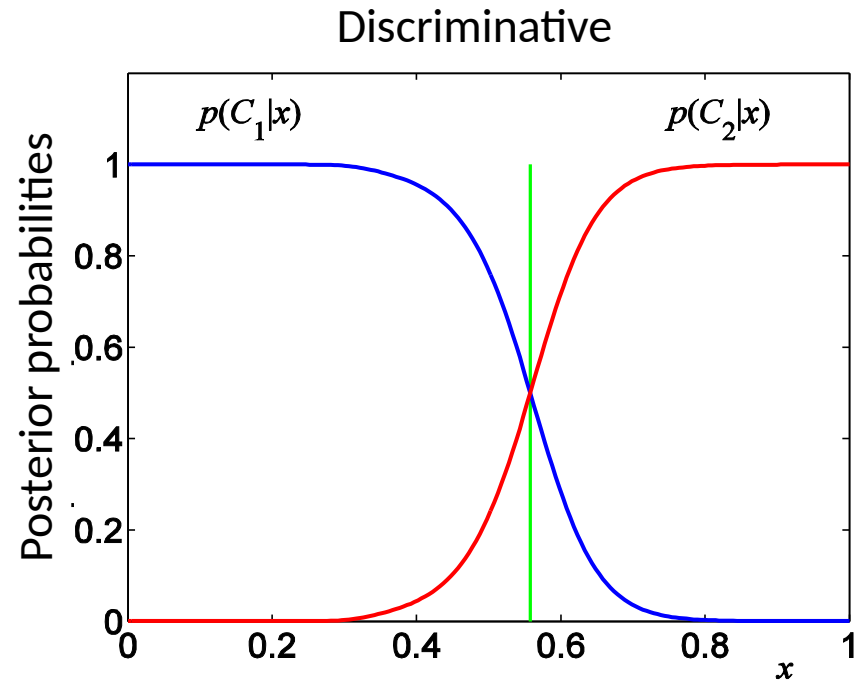
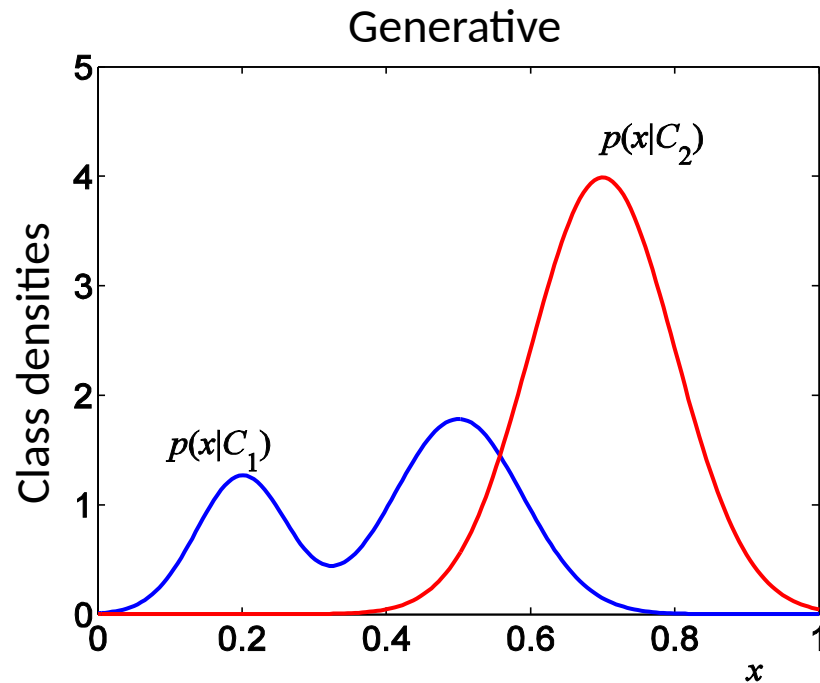
likelihood



posterior (unnormalized)

- Bayesian estimation
 - Goal is to choose the label (skin or \sim skin) that maximizes the posterior \leftrightarrow minimizes probability of misclassification
 - this is called **Maximum A Posteriori (MAP) estimation**

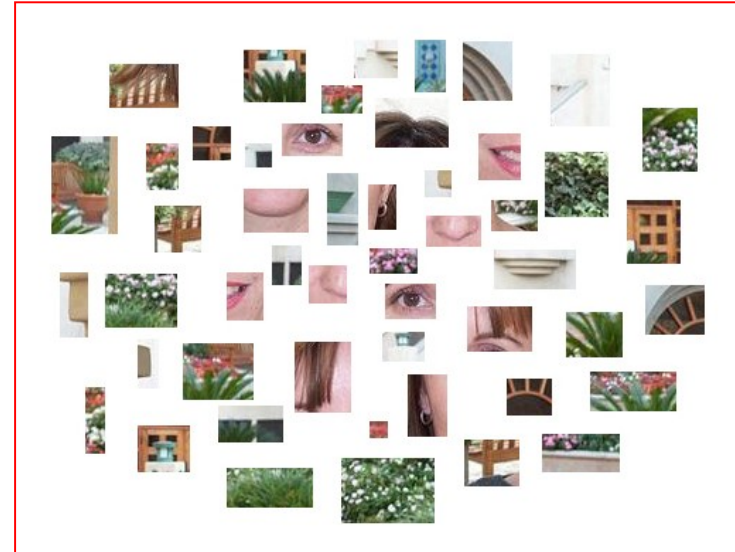
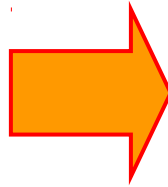
Generative vs. discriminative learning



Steps for statistical recognition

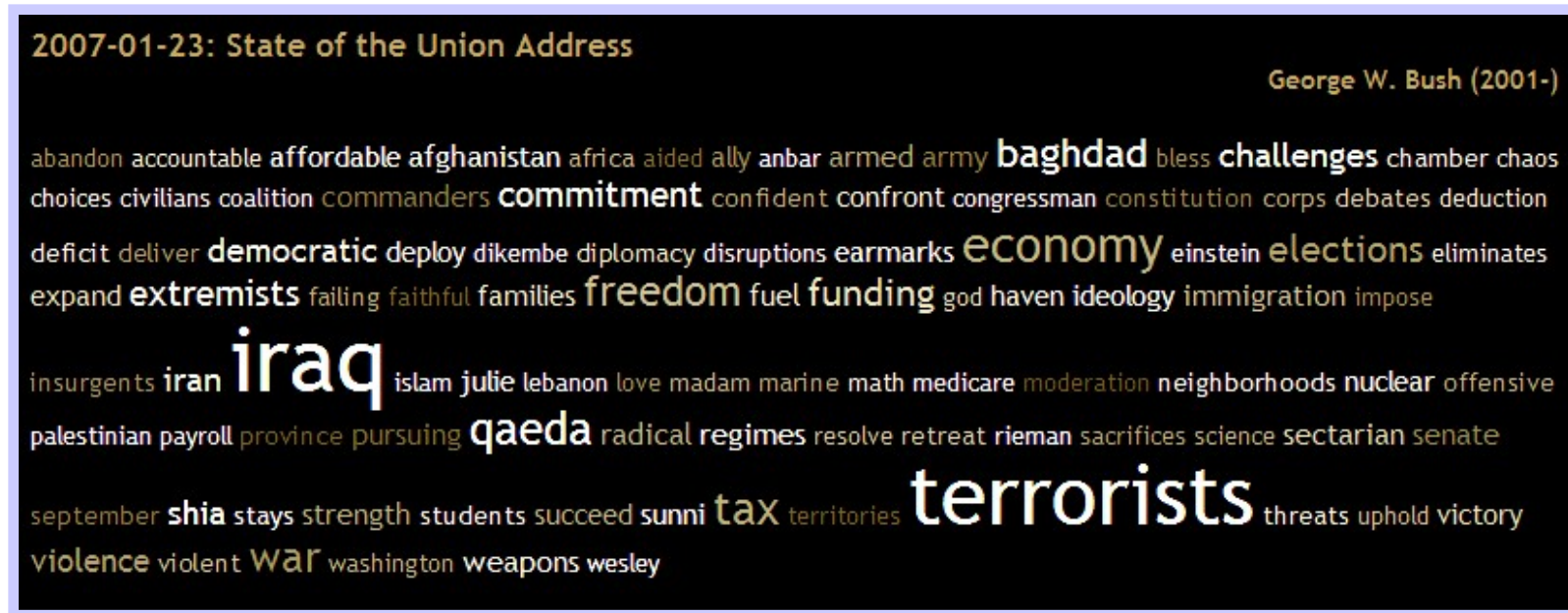
- Representation
 - Specify the model for an object category
 - Bag of features, part-based, global, etc.
- Learning
 - Given a *training set*, find the parameters of the model
 - Generative vs. discriminative
- Recognition
 - Apply the model to a new *test image*

Bag-of-features models



Origin: Bag-of-words models

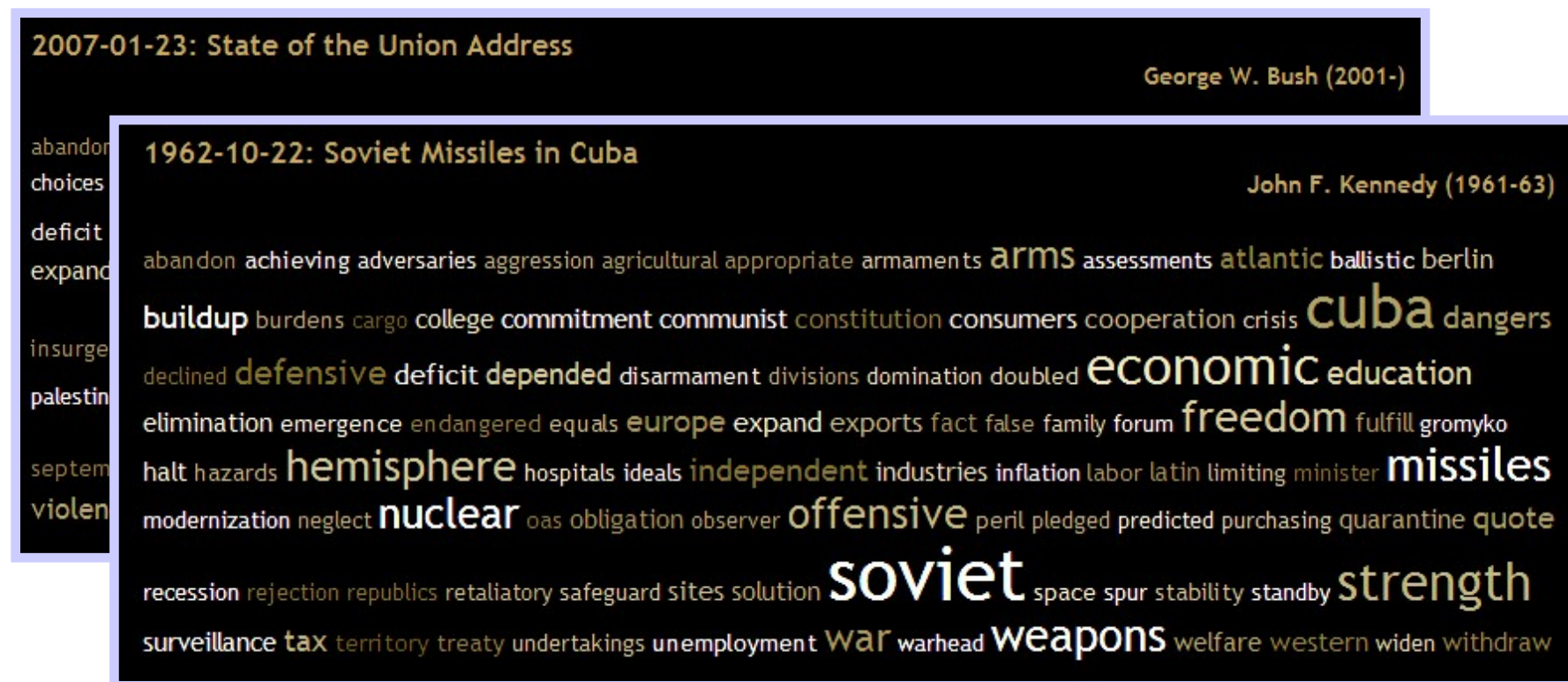
Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



US Presidential Speeches Tag Cloud
<http://chir.ag/projects/preztags/>

Origin: Bag-of-words models

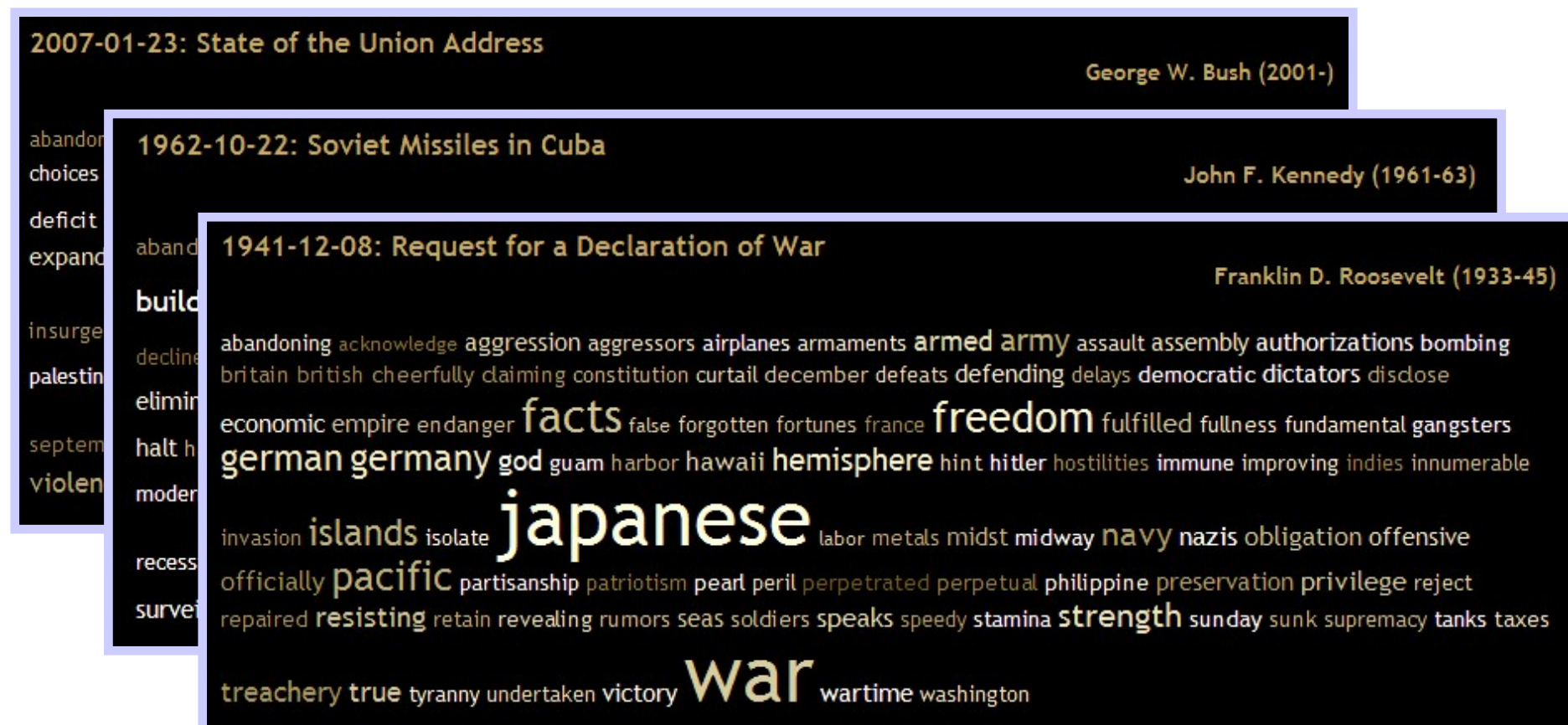
Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



US Presidential Speeches Tag Cloud
<http://chir.ag/projects/preztags/>

Origin: Bag-of-words models

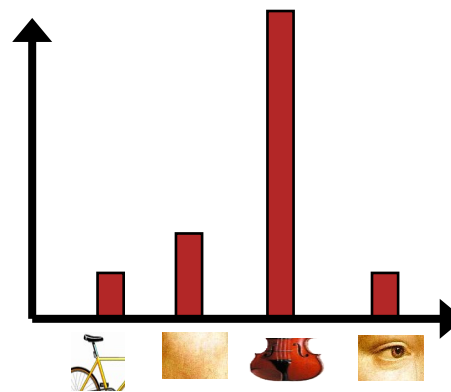
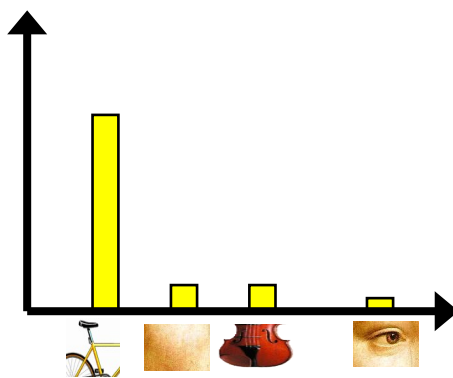
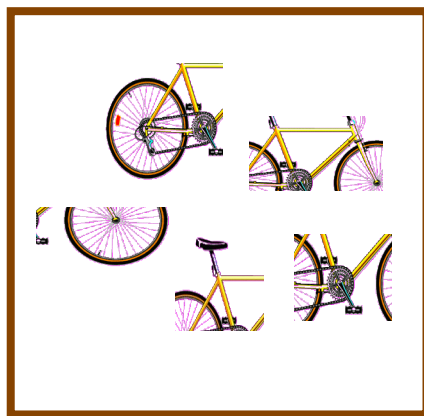
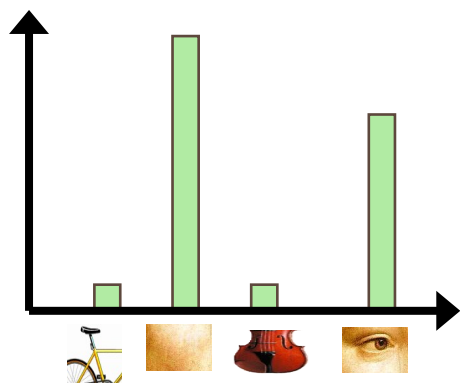
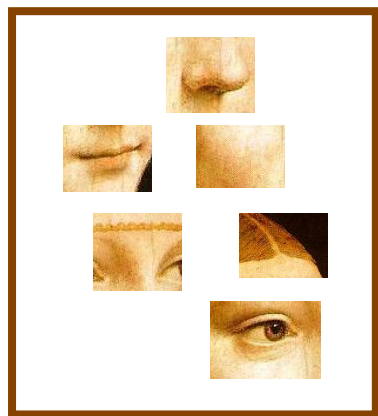
Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



US Presidential Speeches Tag Cloud
<http://chir.ag/projects/preztags/>

Bag-of-features steps

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”

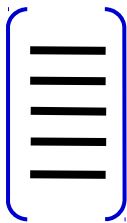


1. Feature extraction

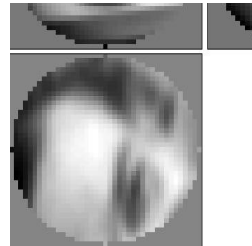
- Regular grid or interest regions



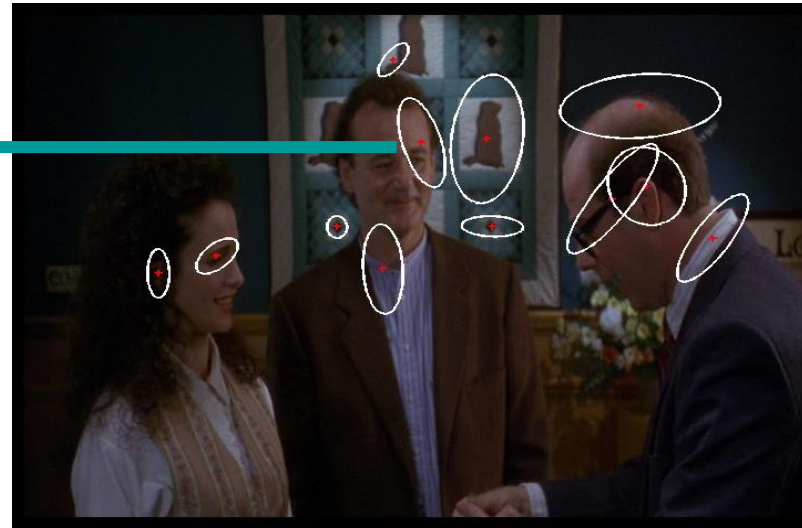
1. Feature extraction



**Compute
descriptor**

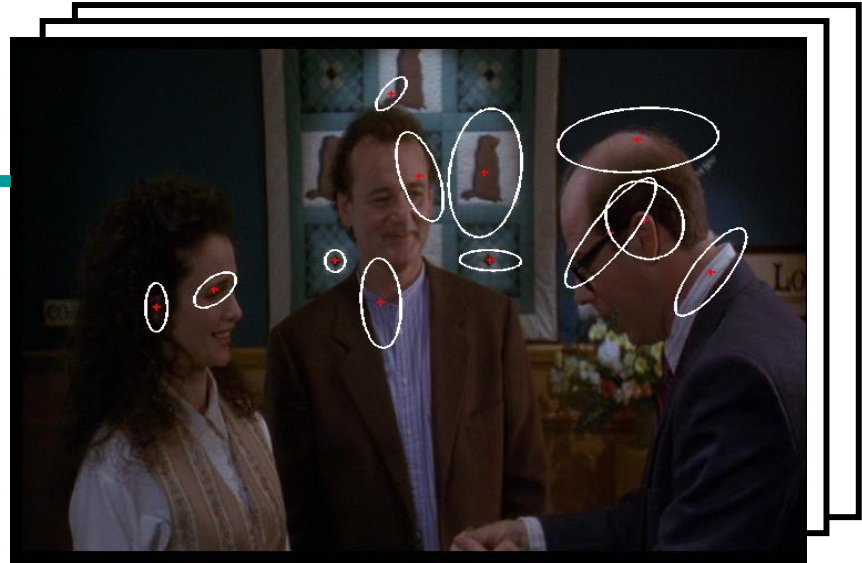
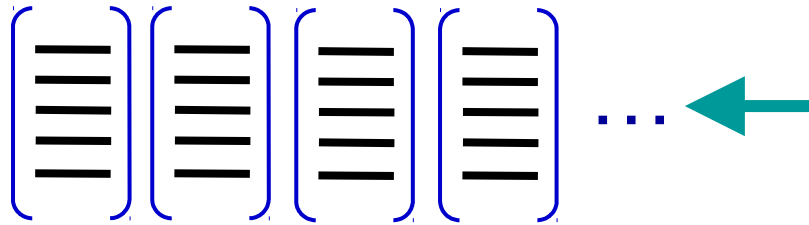


**Normalize
patch**



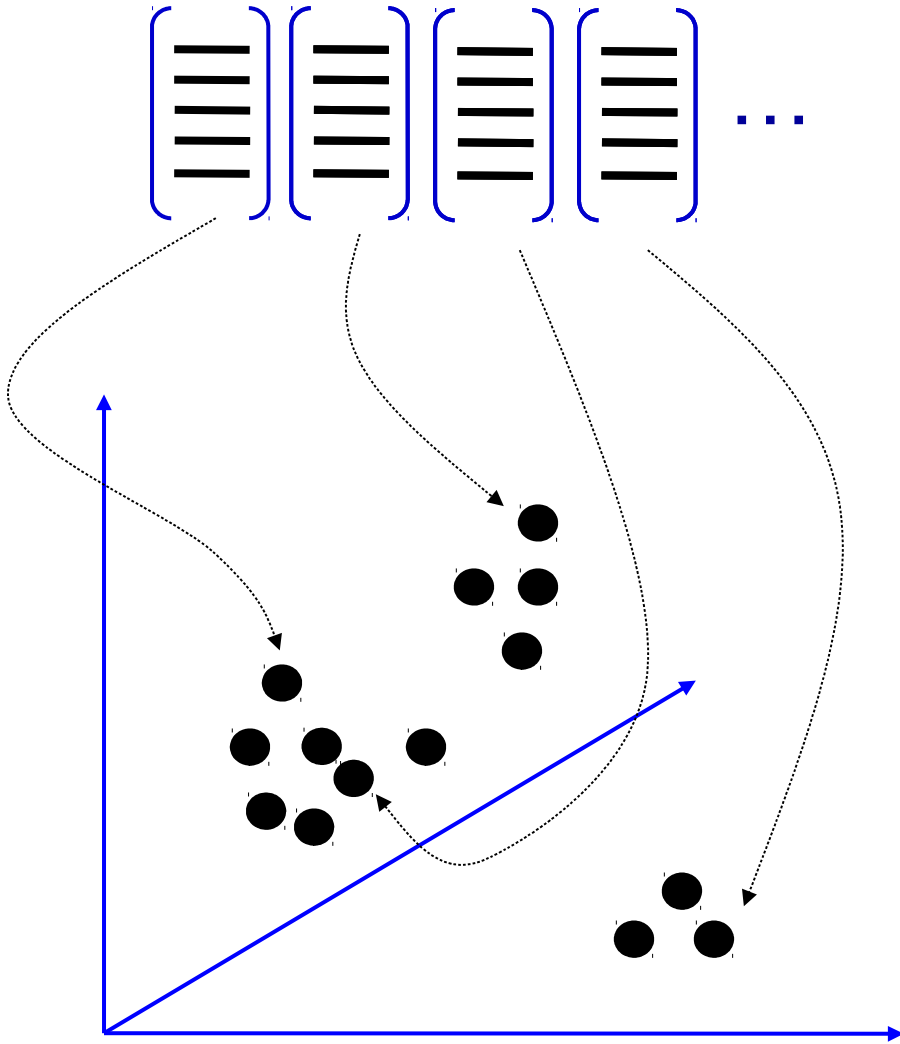
Detect patches

1. Feature extraction



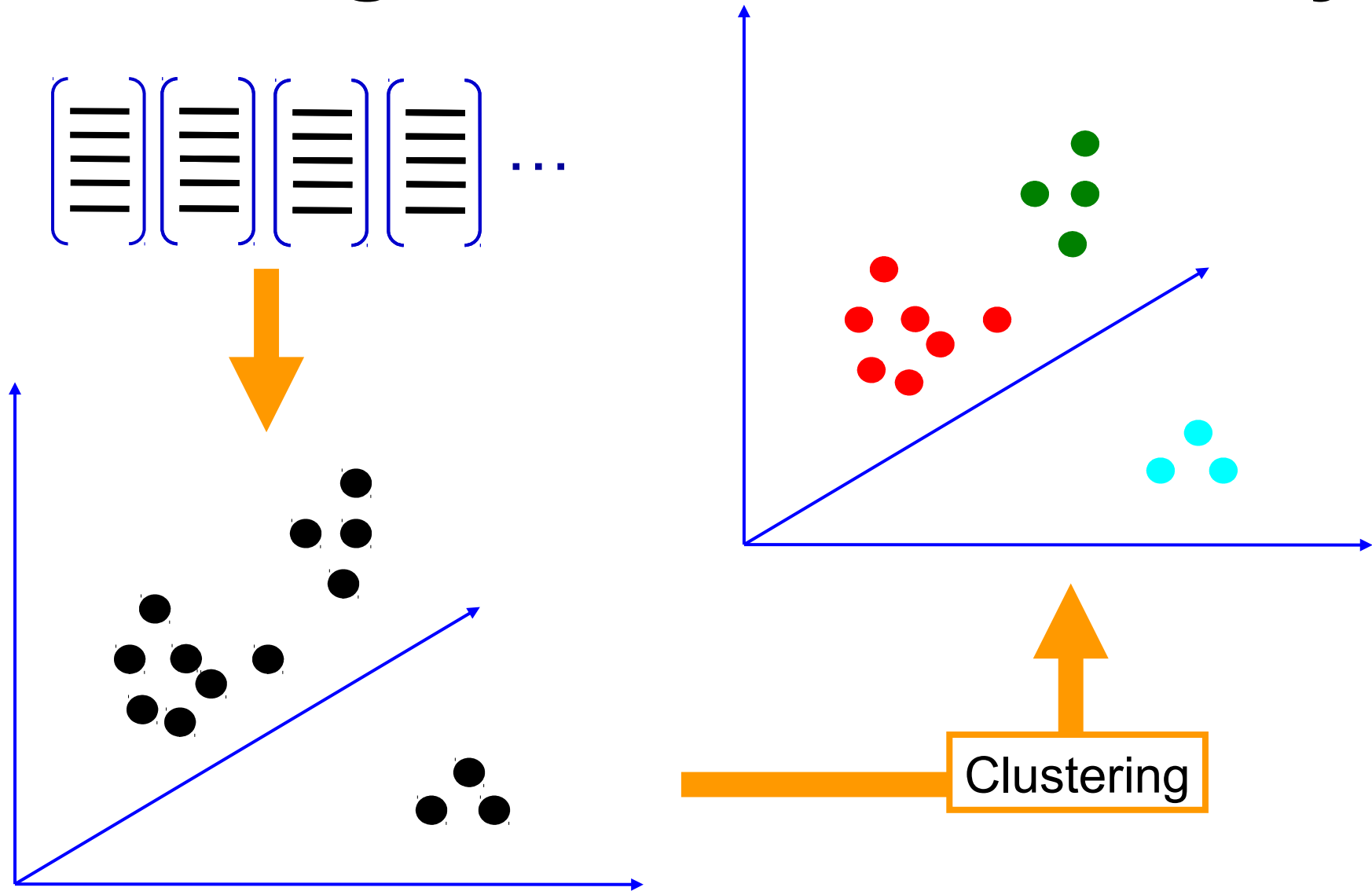
Slide credit: Josef Sivic

2. Learning the visual vocabulary



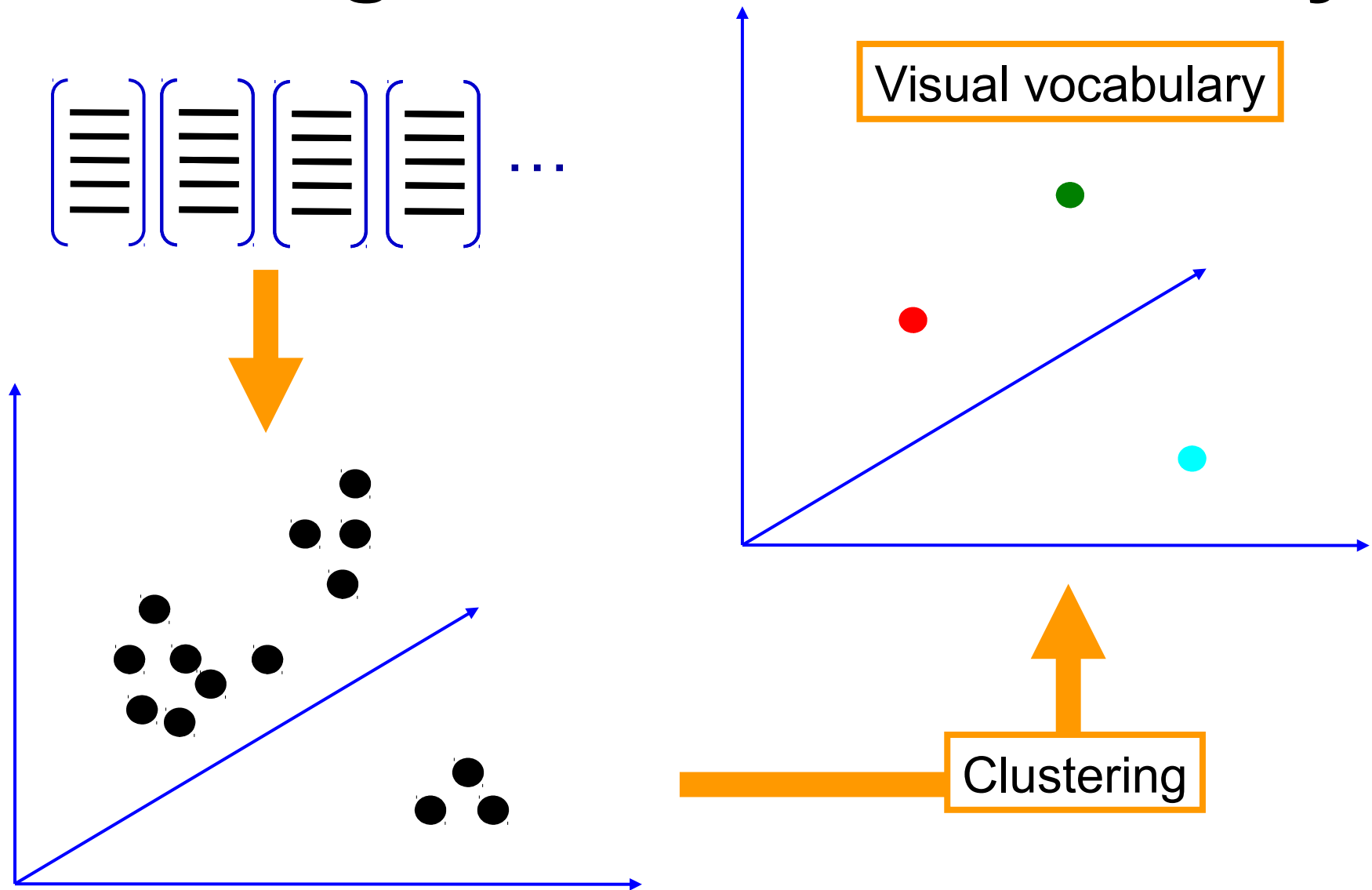
Slide credit: Josef Sivic

2. Learning the visual vocabulary



Slide credit: Josef Sivic

2. Learning the visual vocabulary



Slide credit: Josef Sivic

K-means clustering

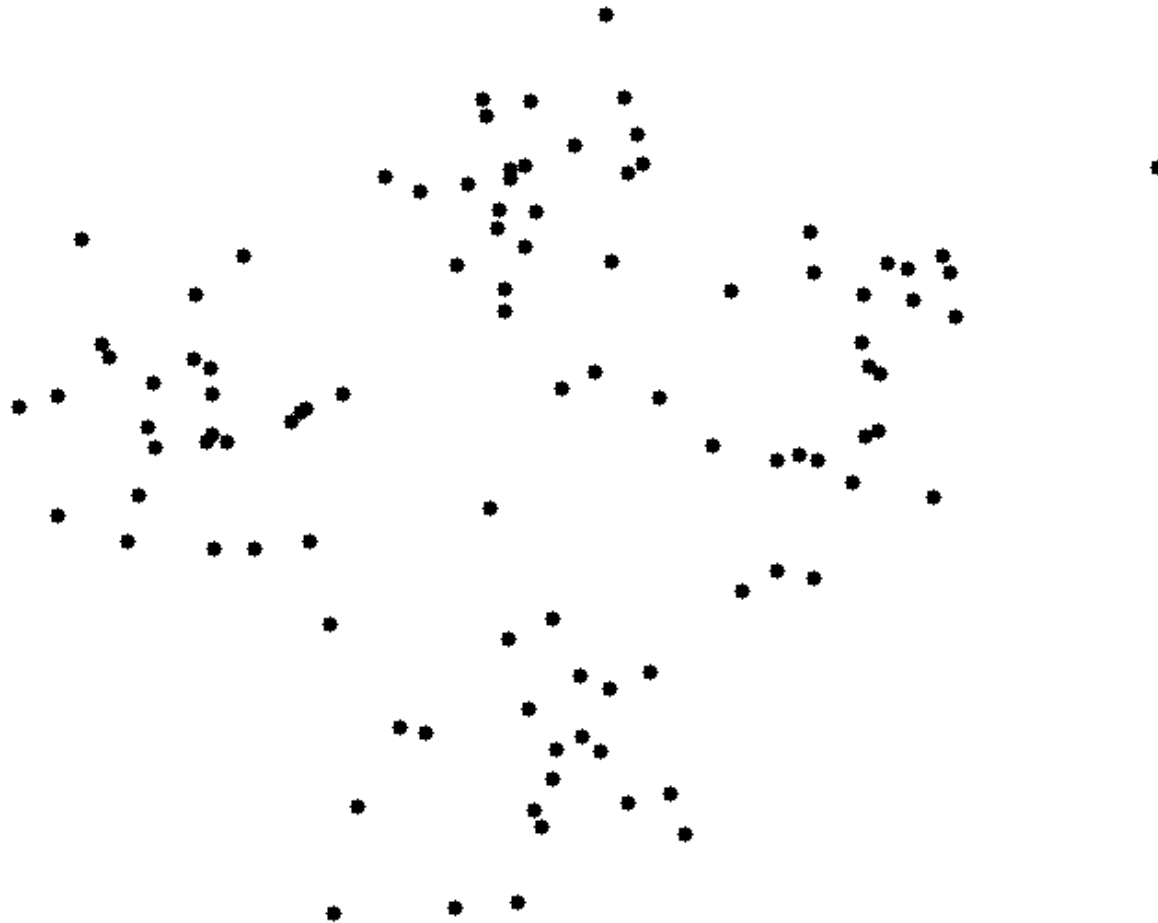
- Want to minimize sum of squared Euclidean distances between points \mathbf{x}_i and their nearest cluster centers \mathbf{m}_k

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\substack{\text{point } i \text{ in} \\ \text{cluster } k}} (\mathbf{x}_i - \mathbf{m}_k)^2$$

Algorithm:

- Randomly initialize K cluster centers
- Iterate until convergence:
 - Assign each data point to the nearest center
 - Recompute each cluster center as the mean of all points assigned to it

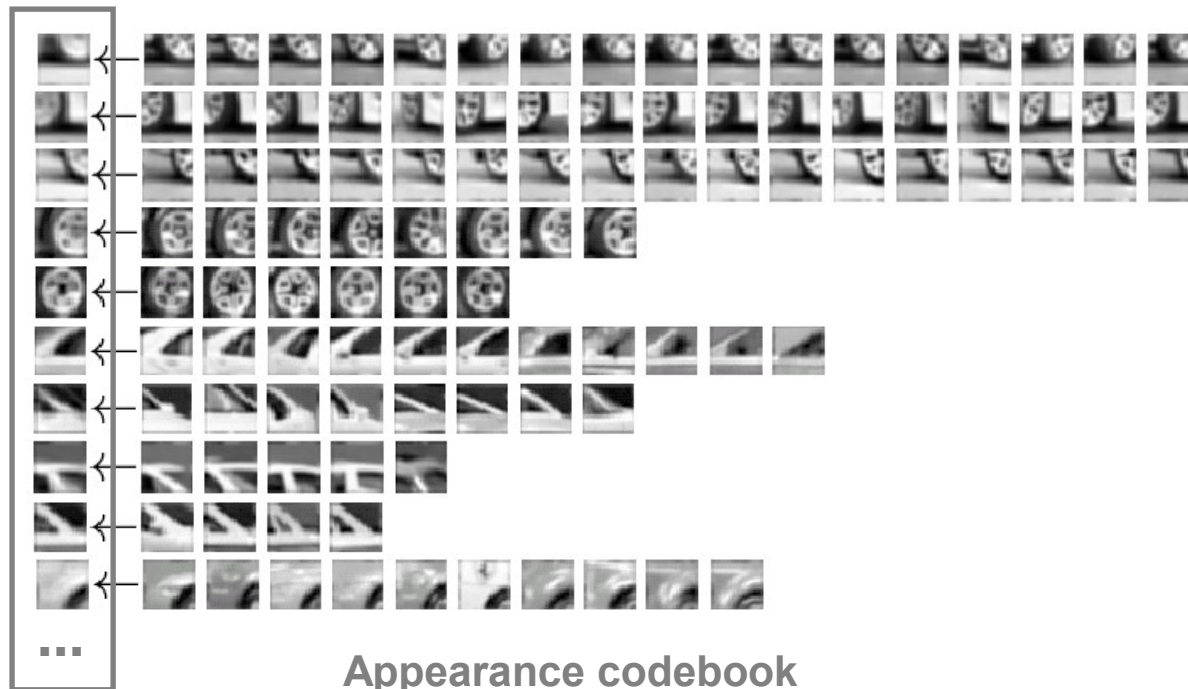
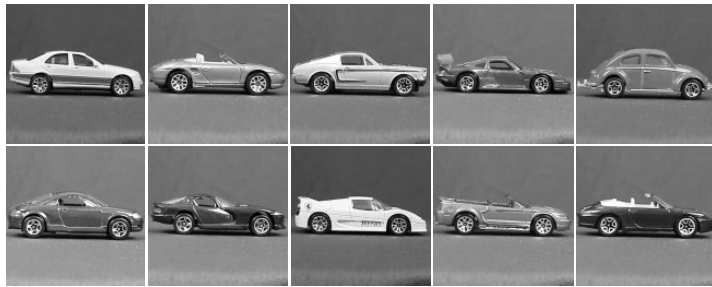
K-means demo



Source: <http://shabal.in/visuals/kmeans/1.html>

Another demo: <http://www.kovan.ceng.metu.edu.tr/~maya/kmeans/>

Example codebook



Source: B. Leibe

Clustering and vector quantization

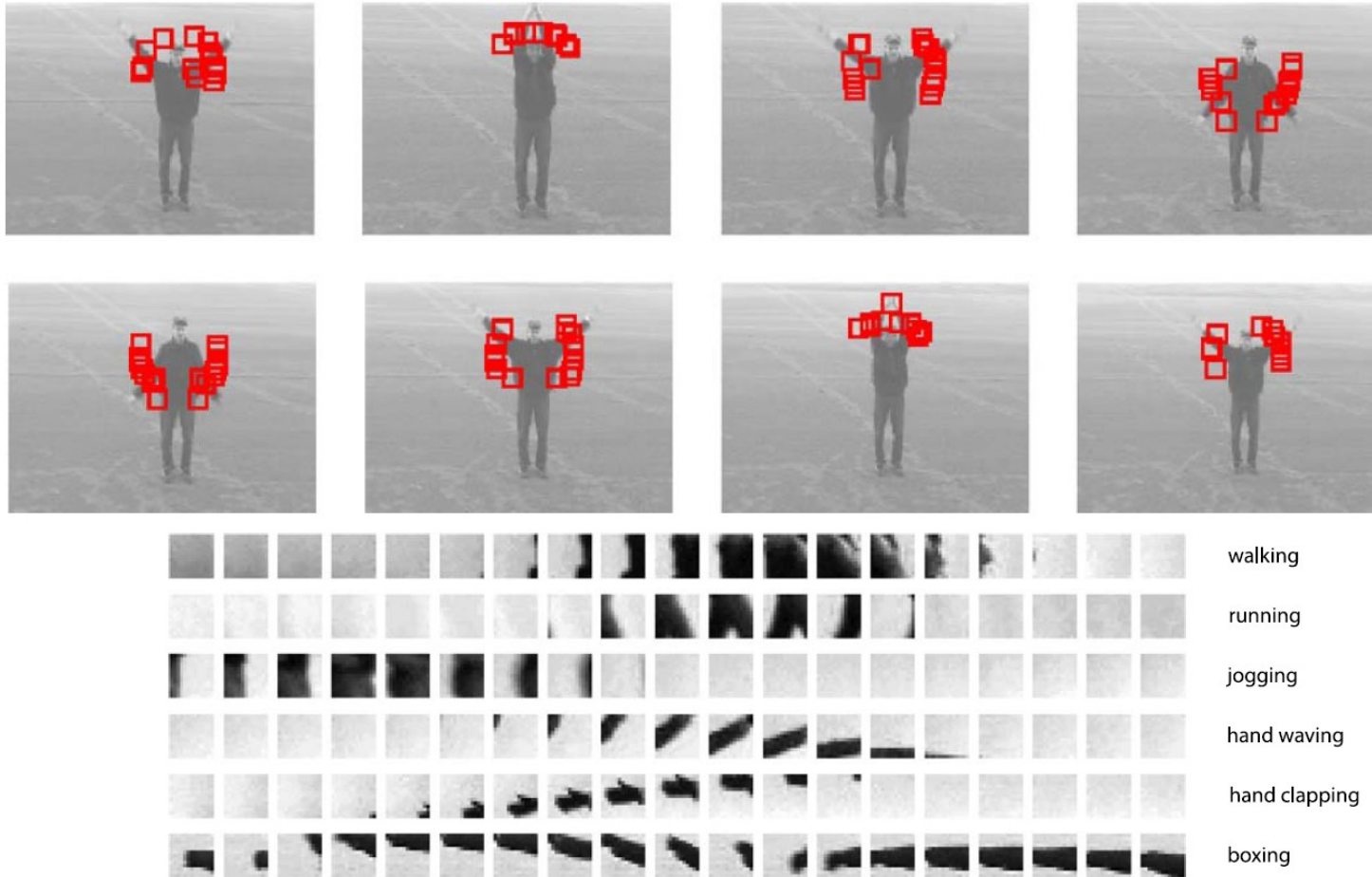
- Clustering is a common method for learning a visual vocabulary or codebook
 - Unsupervised learning process
 - Each cluster center produced by k-means becomes a codevector
 - Codebook can be learned on separate training set
 - Provided the training set is sufficiently representative, the codebook will be “universal”
- The codebook is used for quantizing features
 - A *vector quantizer* takes a feature vector and maps it to the index of the nearest codevector in a codebook
 - Codebook = visual vocabulary
 - Codevector = visual word

Visual vocabularies: Issues

- How to choose vocabulary size?
 - Too small: visual words not representative of all patches
 - Too large: quantization artifacts, overfitting
 - Right size is application-dependent
- Improving efficiency of quantization
 - Vocabulary trees (Nister and Stewenius, 2005)
- Improving vocabulary quality
 - Discriminative/supervised training of codebooks
 - Sparse coding
- More informative bag-of-words representations
 - Fisher Vectors (Perronnin et al., 2007), VLAD (Jegou et al., 2010)
- Incorporating spatial information

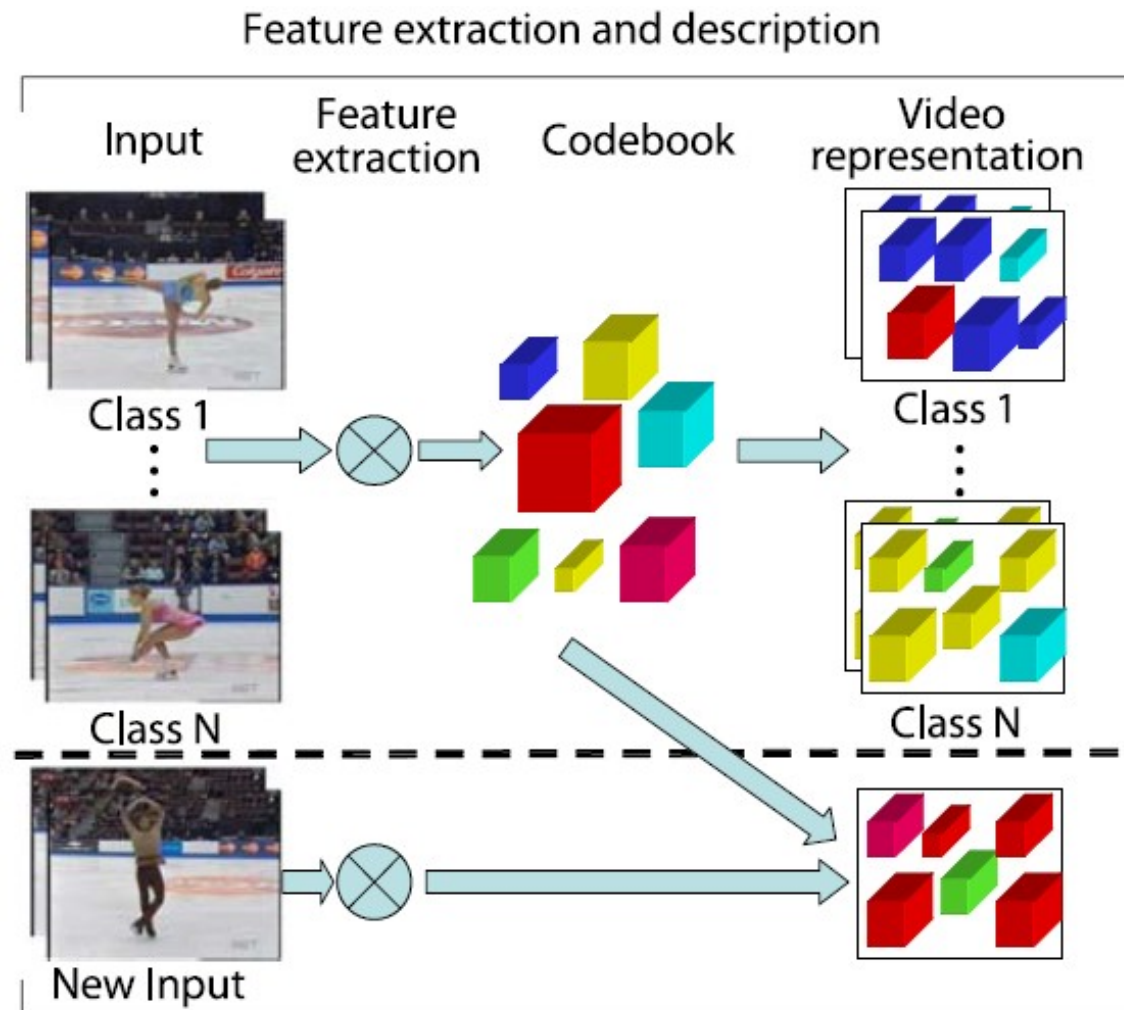
Bags of features for action recognition

Space-time interest points



Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei,
Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words,
IJCV 2008.

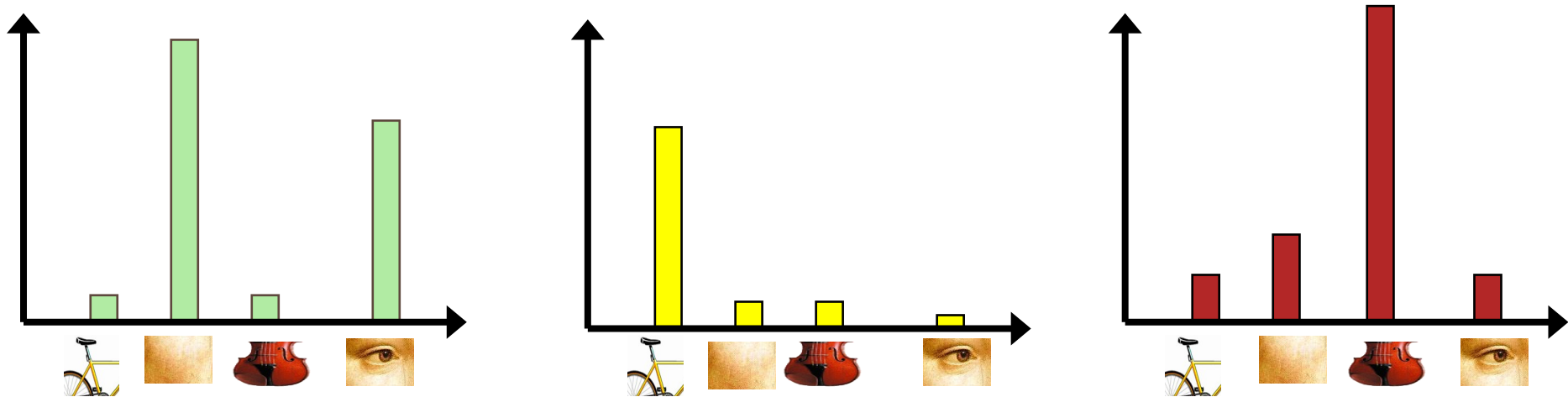
Bags of features for action recognition



Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei,
Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words,
IJCV 2008.

3. classification

- Given the bag-of-features representations of images from different classes, how do we learn a model for distinguishing them?



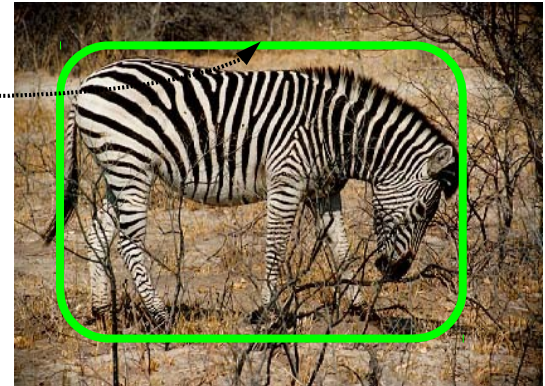
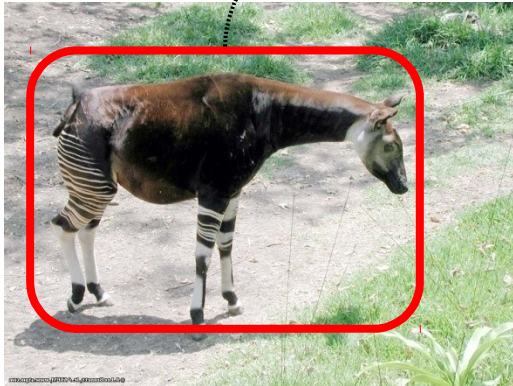
Classifiers

- Learn a decision rule assigning bag-of-features representations of images to different classes

Decision boundary

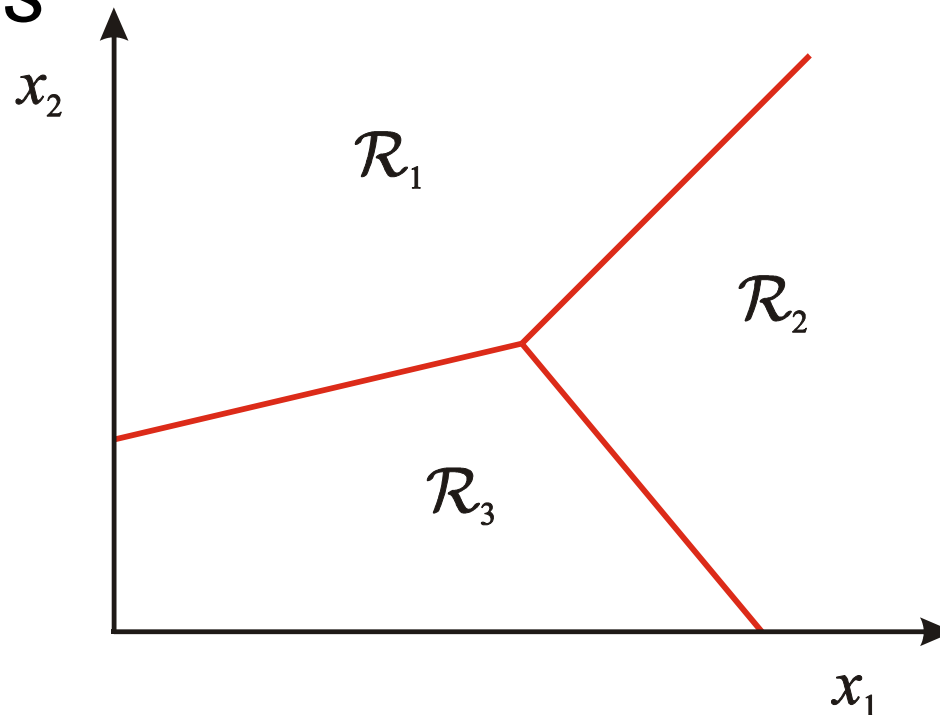
Zebra

Non-zebra



Classification

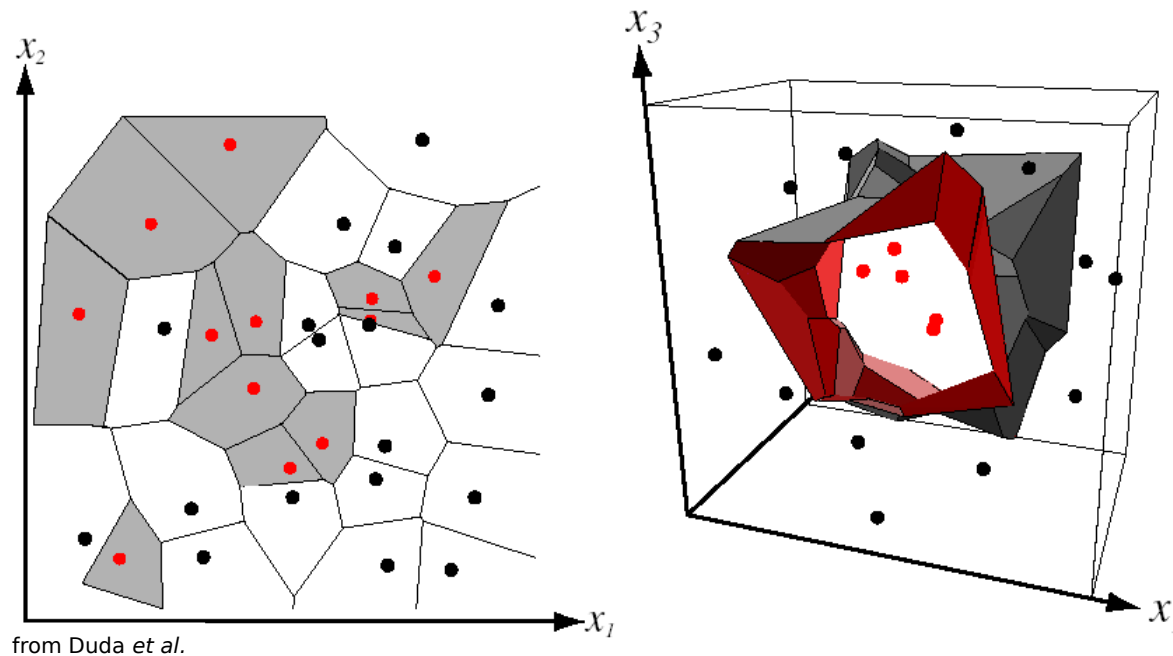
- Assign input vector to one of two or more classes
- Any decision rule divides input space into *decision regions* separated by *decision boundaries*



Nearest Neighbor Classifier

Nearest Neighbor Classifier

- Assign label of nearest training data point to each test data point

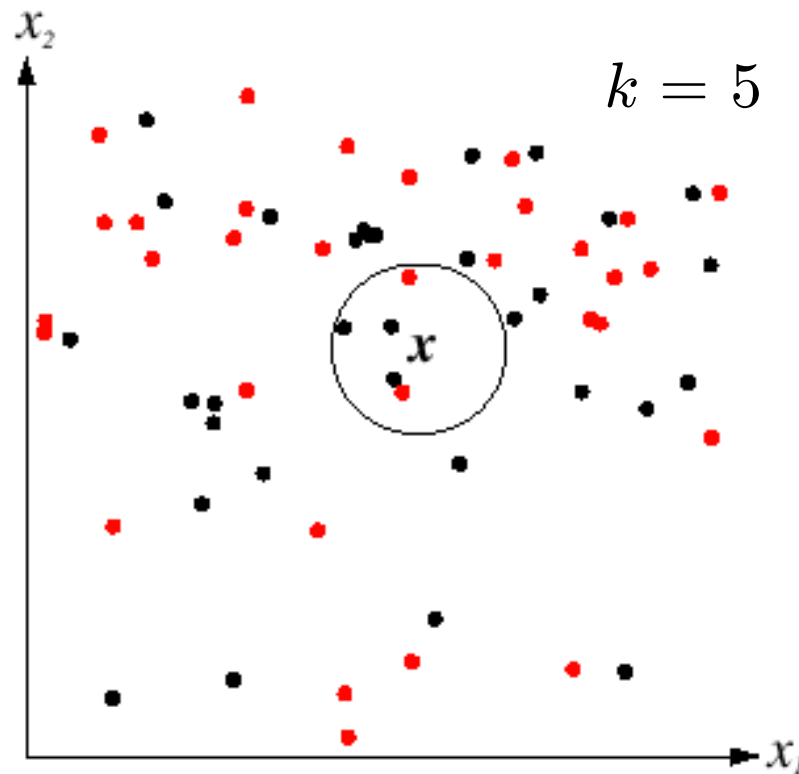


Voronoi partitioning of feature space
for two-category 2D and 3D data

Source: D. Lowe

K-Nearest Neighbors

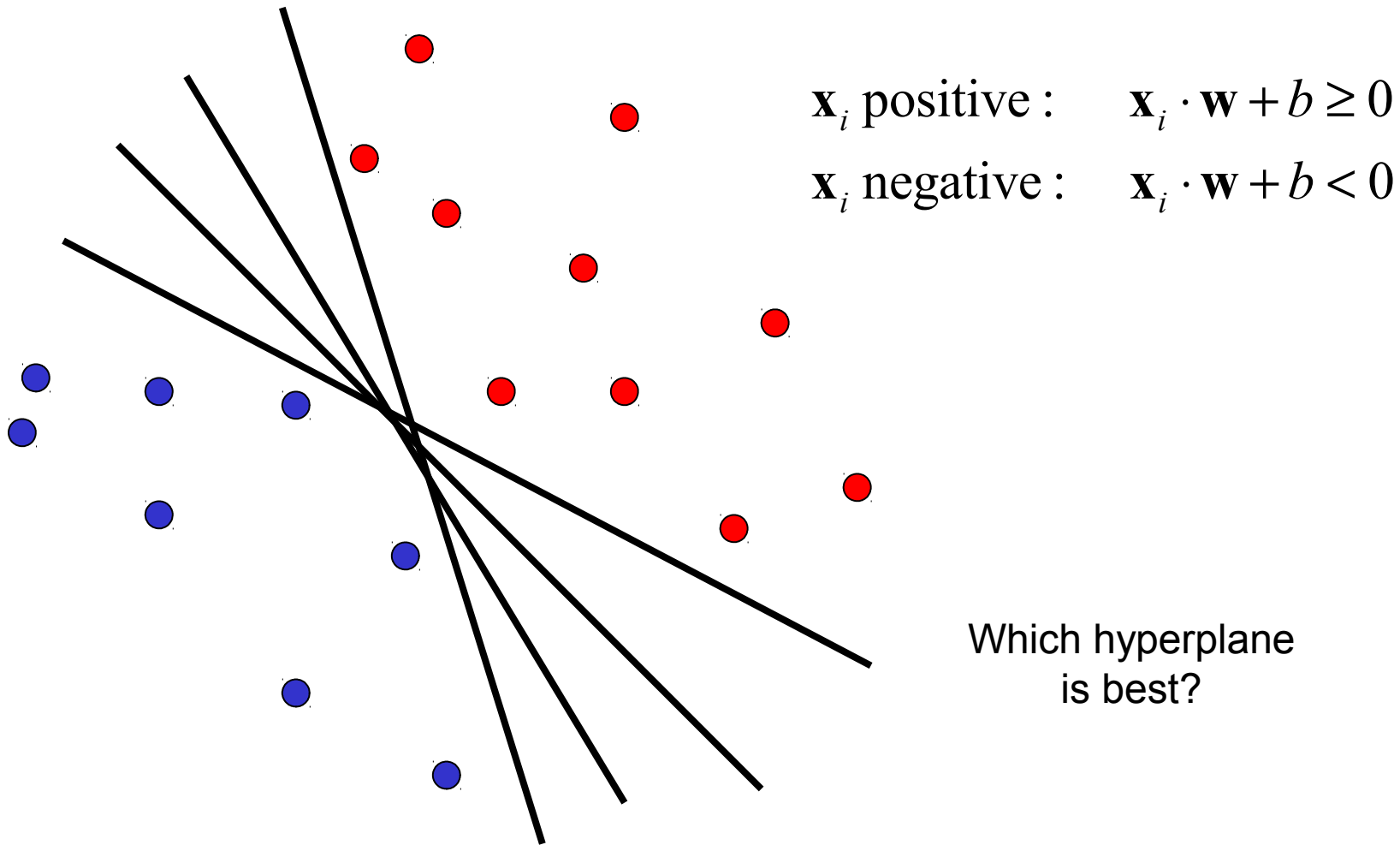
- For a new point, find the k closest points from training data
- Labels of the k points “vote” to classify
- Works well provided there is lots of data and the distance function is good



Source: D. Lowe

Linear classifiers

- Find linear function (*hyperplane*) to separate positive and negative examples

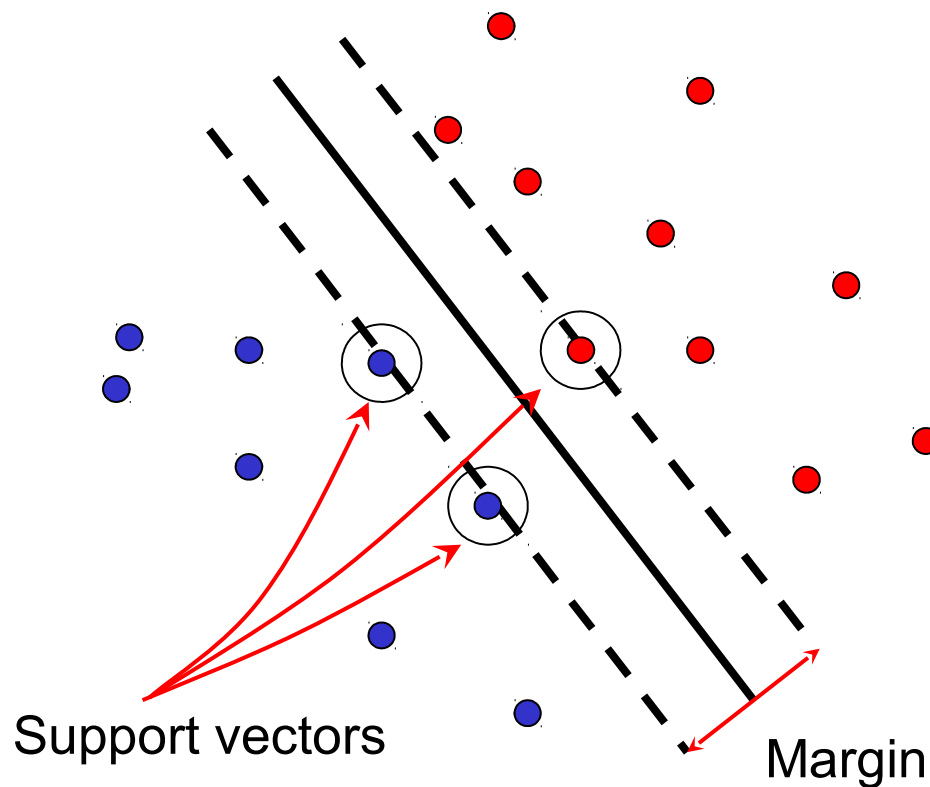


Support vector machines

- Find hyperplane that maximizes the *margin* between the positive and negative examples

Support vector machines

- Find hyperplane that maximizes the *margin* between the positive and negative examples



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

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

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

$$\text{Distance between point and hyperplane:} \quad \frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

$$\text{Therefore, the margin is } 2 / \|\mathbf{w}\|$$

What about multi-class SVMs?

- Unfortunately, there is no “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two-class SVMs
- One vs. others
 - Training: learn an SVM for each class vs. the others
 - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- One vs. one
 - Training: learn an SVM for each pair of classes
 - Testing: each learned SVM “votes” for a class to assign to the test example

SVMs: Pros and cons

- Pros

- Many publicly available SVM packages:
<http://www.kernel-machines.org/software>
- Kernel-based framework is very powerful, flexible
- SVMs work very well in practice, even with very small training sample sizes

- Cons

- No “direct” multi-class SVM, must combine two-class SVMs
- Computation, memory
 - During training time, must compute matrix of kernel values for every pair of examples
 - Learning can take a very long time for large-scale problems
- Unbalanced classes may be problematic

SVM: Practical Advice

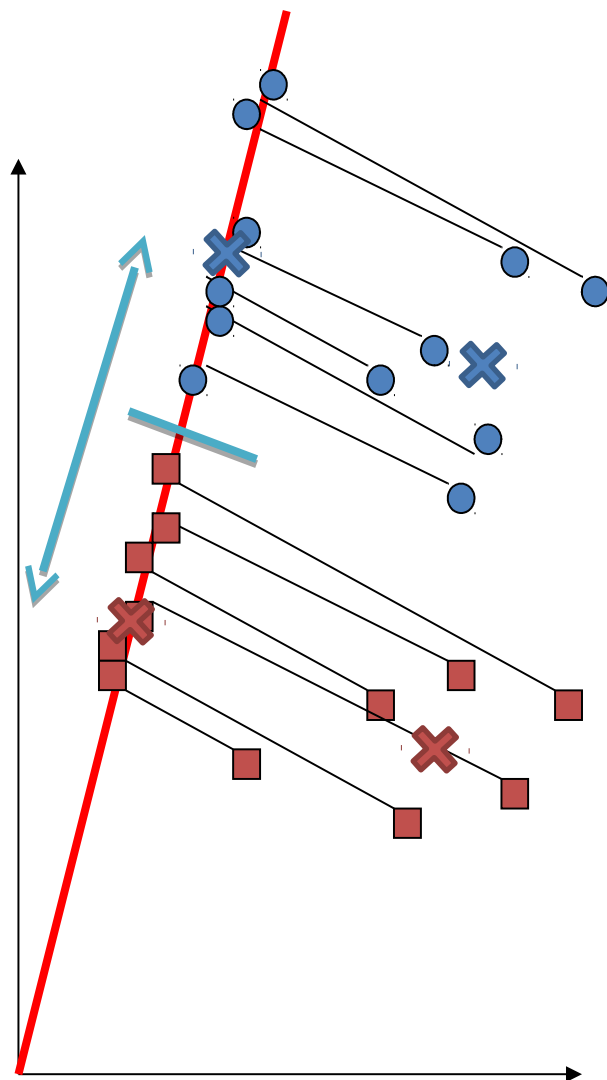
- Avoid MATLAB's implementation
- Use
<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- Or
<http://www.csie.ntu.edu.tw/~cjlin/liblinear/>
when you know you want linear SVM
- Play with C
- Avoid unbalanced classes

LDA

Linear Discriminant Analysis (LDA)

- הרעיון:

1. שימוש ב-LDA למציאת הטלה לממד נמוך
2. בחירת גבול הסיווג על ידי בחירה של היפר מישור מפריד (במקרים רבים, נק. האמצע בין הטלות מרכזי הקבוצות)



LDAs: Pros and cons

- Pros
 - Easy to implement
 - Efficient
 - Works well with unbalanced classes
 - Is often all that's needed!
- Cons
 - No “direct” multi-class LDA, must combine two-class LDAs
 - Linear separator
 - Makes assumptions on the distributions (Gaussians with similar cov.)...but may work well even when violated

Knn classification revisited

Boiman, Shechtman, Irani; “In Defense of Nearest-Neighbor Based Image Classification”, CVPR’08

Non-parametric classifiers

- The most common non-parametric method is kNN
- Base their classification decision directly on the data, and require (almost) no learning/training of parameters.

Advantages over Learning-Based Approaches

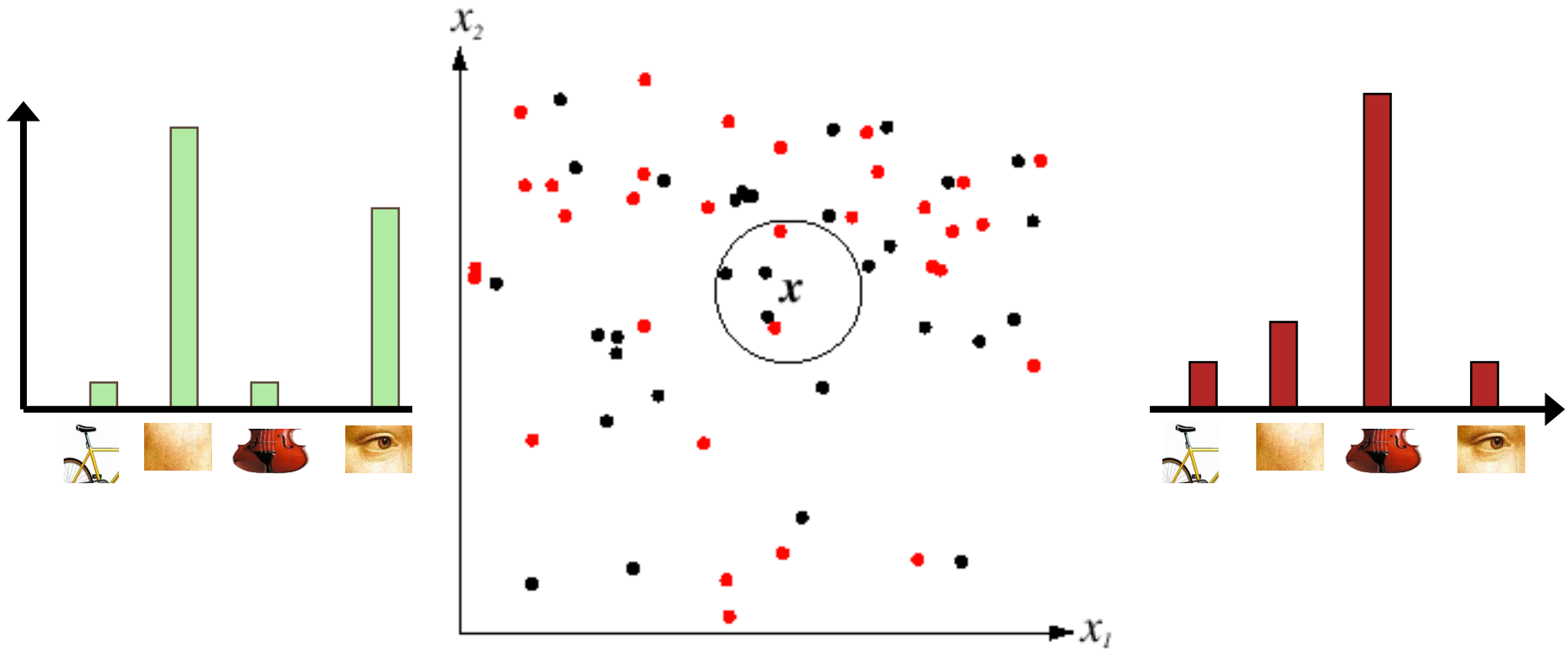
- Naturally handle multiple classes
- No parameters, and so no overfitting of parameters
- Require no training / learning phase

BUT

- Inferior performance compared to learning based approaches (LDA, SVM, AdaBoost, ...)

Why?

- Reason 1: Descriptor quantization
- Reason 2: image-to-image distance



Reason 1: Intuition

- **Highly frequent** descriptors have **low quantization error**, while **rare descriptors** have **high quantization error**.
- The most frequent descriptors in a large database of images (e.g., Caltech-101) comprise of simple edges and corners that appear abundantly in all the classes within the database, and therefore are least informative for classification.
- In contrast, the most informative descriptors for classification are the ones found in one (or few) class, but are rare in other classes.

Reason 1: Example

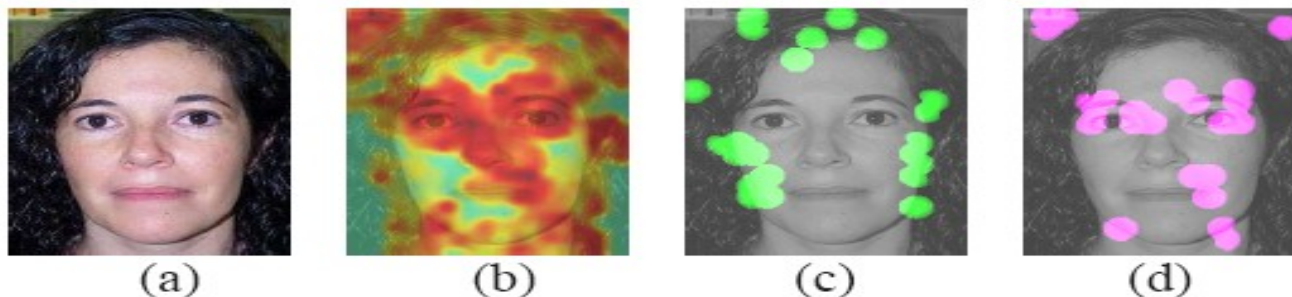
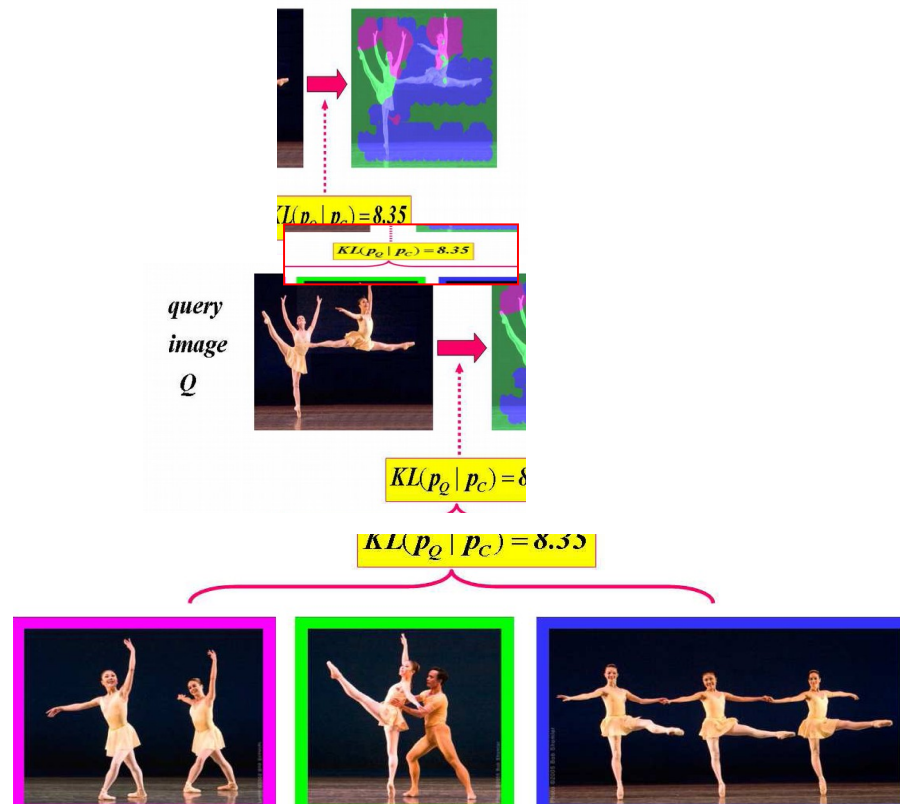


Figure 1. Effects of descriptor quantization – Informative descriptors have low database frequency, leading to high quan-

Reason 2: Image-2-Image

- NN-image classifiers provide good image classification when the query image is similar to one of the labeled images in its class.
- Few example images from a class with large variability \Rightarrow bad classification!

Reason 2: Example



Conclusion

- Data analysis
- PCA
- Classification problem
- Performance evaluation