
INF 4300 – Classification III

Anne Solberg

28.10.15

The agenda today:

- More on estimating classifier accuracy
- Curse of dimensionality and simple feature selection
- kNN-classification
- K-means clustering

28.10.15

INF 4300

1

Confusion matrices

- A matrix with the true class label versus the estimated class labels for each class

Estimated class labels

	Class 1	Class 2	Class 3	Total # of samples
Class 1	80	15	5	100
Class 2	5	140	5	150
Class 3	25	50	125	200
Total	110	205	135	450

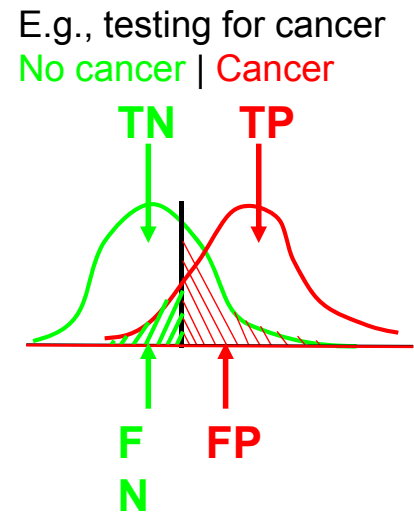
True class labels

INF 4300

2

True / False positives / negatives

- **True positive (TP):**
Patient has cancer
and test result is positive.
- **True negative (TN):**
A healthy patient
and a negative test result.
- **False positive (FP):**
Healthy patient that gets a positive test result.
- **False negative (FN):**
Cancer patient that gets a negative test result.
- *Good to have: TP & TN*
- *Bad to have: FP (but this will probably be detected)*
- *Worst to have: FN (may go un-detected)*

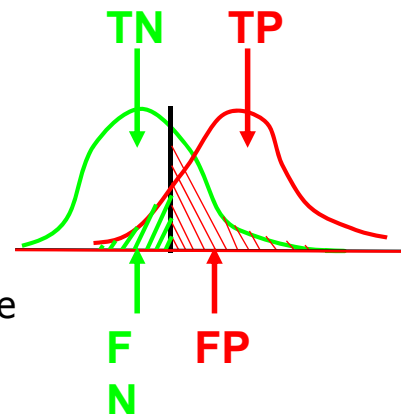


INF 4300

3

Sensitivity and specificity

- **Sensitivity:**
the portion of the data set that tested positive
out of all the positive patients tested:
 - **Sensitivity = $TP / (TP + FN)$**
 - The probability that the test is positive
given that the patient is sick.
 - Higher sensitivity means that
fewer disease cases go undetected.
- **Specificity:**
the portion of the data set that tested negative
out of all the negative patients tested:
 - **Specificity = $TN / (TN + FP)$**
 - The probability that a test is negative
given that the patient is not sick.
 - Higher specificity means that
fewer healthy patients are labeled as sick.



INF 4300

4

Bayes classification with loss functions

- In cases where different classes have different importance (e.g. sick/healthy), we can incorporate this into a Bayesian classifier if we consider the loss.
- Let $\lambda(\alpha_i|\omega_j)$ be the loss if we decide class α_i if the true class is ω_j .
- The risk of deciding class α_i is then: $R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|\mathbf{x})$
- To minimize the overall risk, compute $R(\alpha_i|\mathbf{x})$ for $i=1\dots c$ and choose the class for which $R(\alpha_i|\mathbf{x})$ is minimum.

Outliers and doubt

- In a classification problem, we might want to identify outliers and doubt samples
- We might want an ideal classifier to report
 - ‘this sample is from class l ’ (usual case)
 - ‘this sample is not from any of the classes’ (outlier)
 - ‘this sample is too hard for me’ (doubt/reject)
- The two last cases should lead to a rejection of the sample!

Outliers

- Heuristically defined as "... samples which did not come from the assumed population of samples"
- The outliers can result from some breakdown in preprocessing.
- Outliers can also come from pixels from other classes than the classes in the training data set.
 - Example: K tree species classes, but a few road pixels divide the forest regions.
- One way to deal with outliers is to model them as a separate class, e.g., a gaussian with very large variance, and estimate prior probability from the training data
- Another approach is to decide on some threshold on the a posteriori probability– and if a sample falls below this threshold for all classes, then declare it an outlier.

Doubt samples

- Doubt samples are samples for which the class with the highest probability is not significantly more probable than some of the other classes (e.g. two classes have essentially equal probability).
- Doubt pixels typically occur on the border between two classes ("mixels")
 - Close to the decision boundary the probabilities will be almost equal.
- Classification software can allow the user to specify thresholds for doubt.

The training / test set dilemma

- Ideally we want to maximize the size of both the training and test dataset
- Obviously there is a fixed amount of available data with known labels
- A very simple approach is to separate the dataset in two random subsets
- For small sample sizes we may have to use another strategy: Cross-validation
- This is a good strategy when we have very few "ground truth" samples.
 - Common in medicine where we might have a small number of patients with a certain type of cancer.
 - The cost of obtaining more ground truth data might be so high that we have to do with a small number of ground truth samples.

Crossvalidation / Leave – n - Out

- A very simple (but computationally complex) idea allows us to "fake" a large test set
 - Train the classifier on a set of $N-n$ samples
 - Test the classifier on the n remaining samples
 - Repeat n/N times (dependent on subsampling)
 - Report average performance on the repeated experiments as "test set" error
- An example with leave-1-out and 30 samples:
 - Select one sample to leave out
 - Train on the remaining 29 samples
 - Classify the one sample and store its class label
 - Repeat this 30 times
 - Count the number of misclassifications among the 30 experiments.
- **Leave-n-Out estimation generally overestimates the classification accuracy.**
 - **Feature selection should be performed within the loop, not in advance!!!**
- **Using a training set and a test set of approximately the same size is better.**

The covariance matrix and dimensionality

- Assume we have S classes and a d -dimensional feature vector.
- With a fully multivariate Gaussian model, we must estimate S different mean vectors and S different covariance matrices from training samples.

$\hat{\mu}_s$ has d elements

$\hat{\Sigma}_s$ has $d(d+1)/2$ elements

- Assume that we have M_s training samples from each class
- Given M_s , there is a maximum of the achieved classification performance for a certain value of d
 - increasing n beyond this limit will lead to worse performance.
- Adding more features is not always a good idea!
- Total number of samples given by a rule of thumb: **$M > 10 d S$**
- If we have limited training data, we can use diagonal covariance matrices or regularization

The "curse" of dimensionality

- In practice, the curse means that, for a given sample size, there is a maximum number of features one can add before the classifier starts to degrade.

- For a finite training sample size, the correct classification rate initially increases when adding new features, attains a maximum and then begins to decrease.
- For a high dimensionality, we will need lots of training data to get the best performance.
- => ≈ 10 samples / feature / class.

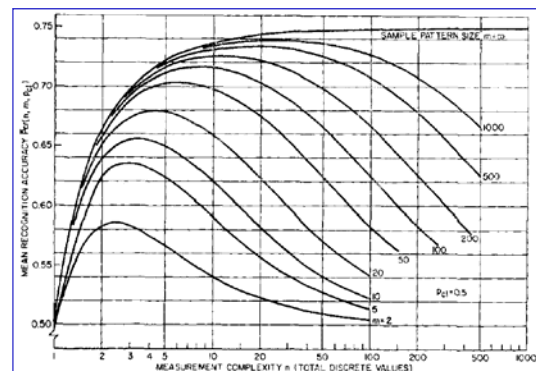


Fig. 3. Finite data set accuracy ($p_{01} = \frac{1}{2}$).

Correct classification rate as function of feature dimensionality, for different amounts of training data. Equal prior probabilities of the two classes is assumed.

Use few, but good features

- To avoid the “curse of dimensionality” we must take care in finding a set of relatively few features.
- A good feature has high within-class homogeneity, and should ideally have large between-class separation.
- In practise, one feature is not enough to separate all classes, but a good feature should:
 - separate some of the classes well
 - Isolate one class from the others.
- If two features look very similar (or have high correlation), they are often redundant and we should use only one of them.
- Class separation can be studied by:
 - Visual inspection of the feature image overlaid the training mask
 - Scatter plots
- Evaluating features as done by training can be difficult to do automatically, so manual interaction is normally required.

How do we beat the “curse of dimensionality”?

- Use regularized estimates for the Gaussian case
 - Use diagonal covariance matrices
 - Apply regularized covariance estimation
- Generate few, but informative features
 - Careful feature design given the application
- Reducing the dimensionality
 - Feature selection – select a subset of the original features (more in INF5300)
 - Feature transforms – compute a new subset of features based on a linear combination of all features (INF 5300)
 - Example 1: Principal component transform
 - Unsupervised, finds the combination that maximized the variance in the data.
 - Example 2: Fisher’s linear discriminant
 - Supervised, finds the combination that maximizes the distance between the classes.

Regularized covariance matrix estimation

- Let the covariance matrix be a weighted combination of a class-specific covariance matrix Σ_k and a common covariance matrix Σ (estimated from training samples for all classes) :

$$\Sigma_k(\alpha) = \frac{(1-\alpha)n_k\Sigma_k + \alpha n\Sigma}{(1-\alpha)n_k + \alpha n}$$

where $0 \leq \alpha \leq 1$ must be determined, and n_k and n is the number of training samples for class k and overall.

- Alternatively:

$$\Sigma_k(\beta) = (1-\beta)\Sigma_k + \beta I$$

where the parameter $0 \leq \beta \leq 1$ must be determined.

- The effect of these are that we can use a quadratic classifier even if we have little training data/ill-conditioned Σ_k
- We still have to be able to compute Σ_k , but the only the regularized/more robust $\Sigma_k(\alpha)$ or $\Sigma_k(\beta)$ must be inverted.

Exhaustive feature selection

- If – for some reason – you know that you will use d out of D available features, an exhaustive search will involve a number of combinations to test:

$$n = \frac{D!}{(D-d)! d!}$$

- If we want to perform an exhaustive search through D features for the optimal subset of the $d \leq m$ “best features”, the number of combinations to test is

$$n = \sum_{d=1}^m \frac{D!}{(D-d)! d!}$$

- Impractical even for a moderate number of features!

$$d \leq 5, D = 100 \Rightarrow n = 79.374.995$$

Suboptimal feature selection

- Select **the best single features** based on some quality criteria, e.g., estimated correct classification rate.
 - A combination of the best single features will often imply correlated features and will therefore be suboptimal .
- More in INF 5300
- **“Sequential forward selection”** implies that when a feature is selected or removed, this decision is final.
- **“Stepwise forward-backward selection”** overcomes this.
 - A special case of the **“add - a, remove - r algorithm”**.
- Improved into **“floating search”** by making the number of forward and backward search steps data dependent.
 - “Adaptive floating search”
 - “Oscillating search”.

INF 4300

17

Distance measures used in feature selection

- In feature selection, each feature combination must be ranked based on a criterion function.
- Criteria functions can either be distances between classes, or the classification accuracy on a validation test set.
- If the criterion is based on e.g. the mean values/covariance matrices for the training data, distance computation is fast.
- Better performance at the cost of higher computation time is found when the classification accuracy on a validation data set (different from training and testing) is used as criterion for ranking features.
 - This will be slower as classification of the validation data needs to be done for every combination of features.

INF 5300

18

Distance measures between classes

- How do we compute the distance between two classes:
 - Distance between the closest two points?
 - Maximum distance between two points?
 - Distance between the class means?
 - Average distance between points in the two classes?
 - Which distance measure?
 - Euclidean distance or Mahalanobis distance?
- Distance between K classes:
 - How do we generalize to more than two classes?
 - Average distance between the classes?
 - Smallest distance between a pair of classes?

Class separability measures

- How do we get an indication of the separability between two classes?
 - Euclidean distance between class means $|\mu_r - \mu_s|$
 - Bhattacharyya distance
 - Can be defined for different distributions
 - For Gaussian data, it is

$$B = \frac{1}{8} (\mu_r - \mu_s)^T \left(\frac{\Sigma_r + \Sigma_s}{2} \right)^{-1} (\mu_r - \mu_s) + \frac{1}{2} \ln \left| \frac{\frac{1}{2} (\Sigma_r + \Sigma_s)}{\sqrt{|\Sigma_r| |\Sigma_s|}} \right|$$

- Mahalanobis distance between two classes:

$$\Delta = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)$$

$$\Sigma = N_1 \Sigma_1 + N_2 \Sigma_2$$

Examples of feature selection from INF 5300 - Method 1 - Individual feature selection

- Each feature is treated individually (no correlation/covariance between features is considered)
- Select a criteria, e.g. a distance measure
- Rank the feature according to the value of the criteria $C(k)$
- Select the set of features with the best individual criteria value
- Multiclass situations:
 - Average class separability or
 - $C(k) = \min \text{distance}(i,j)$ - worst case ← Often used
- Advantage with individual selection: computation time
- Disadvantage: no correlation is utilized.

Method 2 - Sequential backward selection

- Select l features out of d
- Example: 4 features x_1, x_2, x_3, x_4
- Choose a criterion C and compute it for the vector $[x_1, x_2, x_3, x_4]^T$
- Eliminate one feature at a time by computing $[x_1, x_2, x_3]^T$, $[x_1, x_2, x_4]^T$, $[x_1, x_3, x_4]^T$ and $[x_2, x_3, x_4]^T$
- Select the best combination, say $[x_1, x_2, x_3]^T$.
- From the selected 3-dimensional feature vector eliminate one more feature, and evaluate the criterion for $[x_1, x_2]^T$, $[x_1, x_3]^T$, $[x_2, x_3]^T$ and select the one with the best value.
- Number of combinations searched:
 $1 + 1/2((d+1)d - l(l+1))$

Method 3: Sequential forward selection

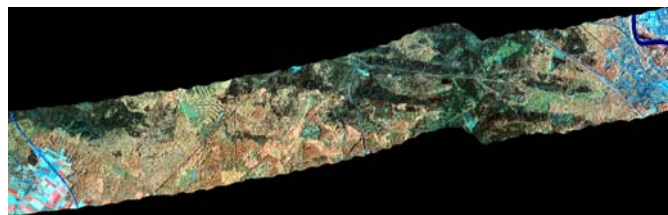
- Compute the criterion value for each feature. Select the feature with the best value, say x_1 .
- Form all possible combinations of features x_1 (the winner at the previous step) and a new feature, e.g. $[x_1, x_2]^T$, $[x_1, x_3]^T$, $[x_1, x_4]^T$, etc. Compute the criterion and select the best one, say $[x_1, x_3]^T$.
- Continue with adding a new feature.
- Number of combinations searched: $d - (d-1)/2$.
 - Backwards selection is faster if l is closer to d than to 1.

INF 4300

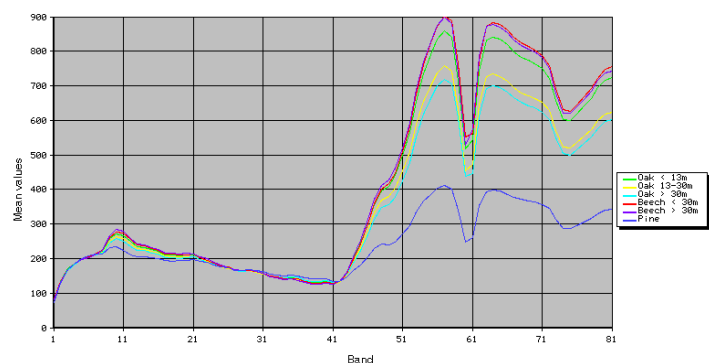
23

Hyperspectral image example

- A hyperspectral image from France
- 81 features/spectral bands
- 6 classes (tree species)
- μ has 81 parameters to compute for each class
- Σ has $81 \times 80 / 2 = 3240$ parameters for each class.
- 1000 training samples for each class.
- Test set: 1000-2000 samples for each class.



3 of the 81 bands shown as RGB image



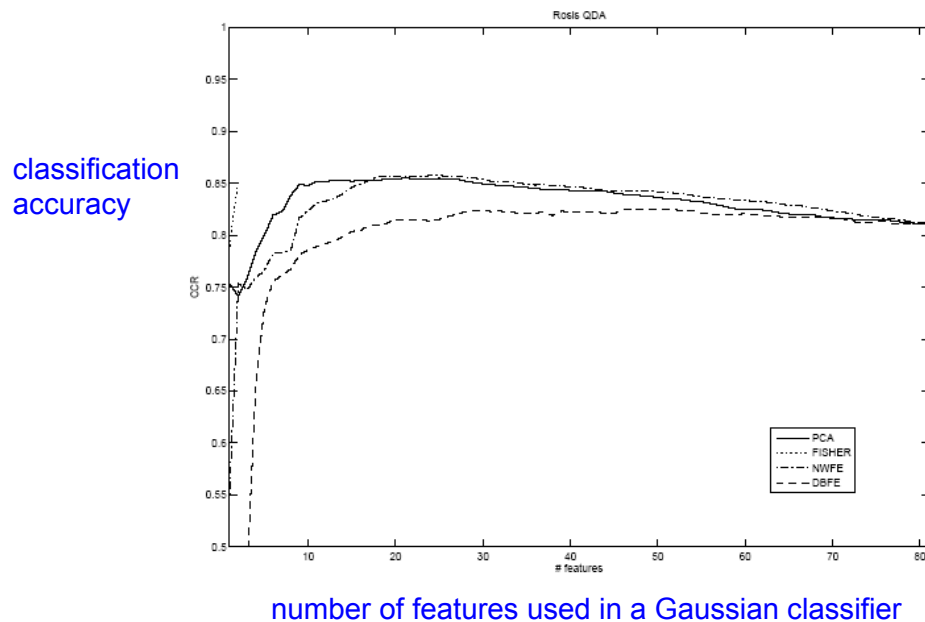
Plots of the 81 mean values for each class

INF 4300

24

Hyperspectral example

classification accuracy vs. nof. features on test set



Each curve shows a different dimensionality reduction method.

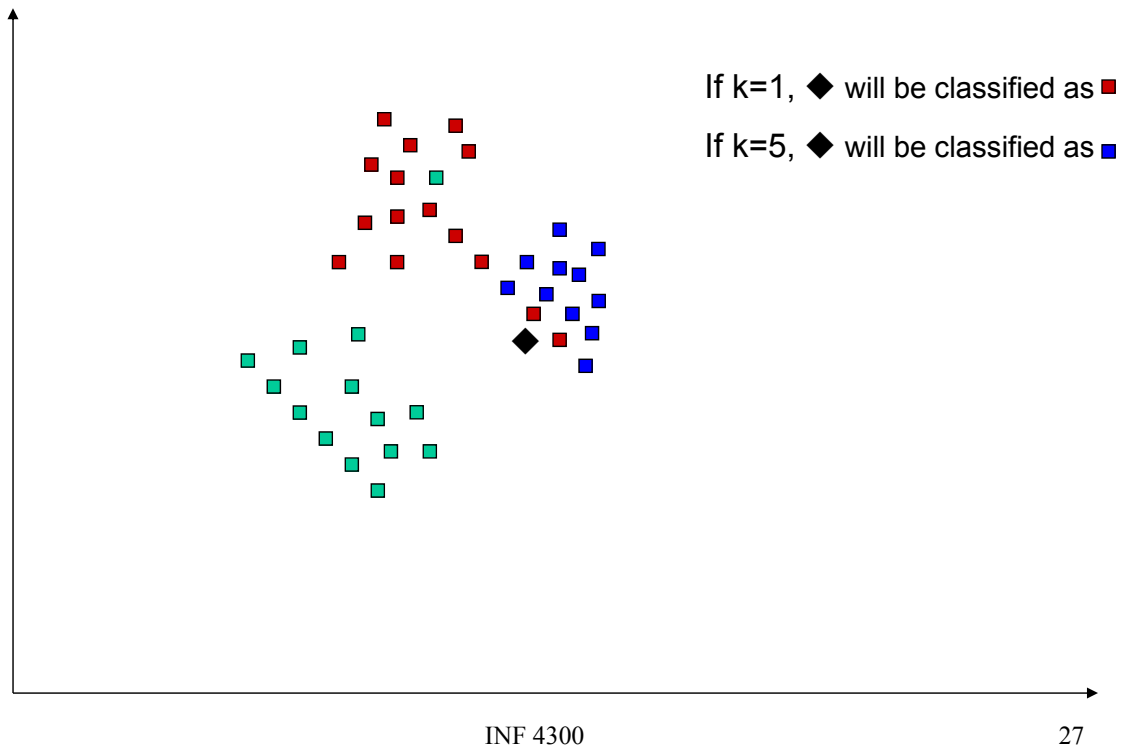
Note that as we include more features, the classification accuracy first increases, then it starts to decrease.

Curse of dimensionality!

k-Nearest-Neighbor classification

- A very simple classifier.
- Classification of a new sample x_i is done as follows:
 - Out of N training vectors, identify the k nearest neighbors (measured by Euclidean distance) in the training set, irrespectively of the class label.
 - Out of these k samples, identify the number of vectors k_j that belong to class ω_i , $i:1,2,\dots,M$ (if we have M classes)
 - Assign x_i to the class ω_i with the maximum number of k_j samples.
- k should be odd, and must be selected a priori.

kNN-example



About kNN-classification

- If $k=1$ (1NN-classification), each sample is assigned to the same class as the closest sample in the training data set.
- If the number of training samples is very high, this can be a good rule.
- If $k \rightarrow \infty$, this is theoretically a very good classifier.
- This classifier involves no "training time", but the time needed to classify one pattern x_i will depend on the number of training samples, as the distance to all points in the training set must be computed.
- "Practical" values for k : $3 \leq k \leq 9$
- *Classification performance should **always** be computed on the test data set.*

Supervised or unsupervised classification

- Supervised classification
 - Classify each object or pixel into a set of k known classes
 - Class parameters are estimated using a set of **training samples** from each class.
- Unsupervised classification
 - Partition the feature space into a set of k clusters
 - k is not known and must be estimated (difficult)
- In both cases, classification is based on the value of the set of n features x_1, \dots, x_n .
- The object is classified to the class which has the highest posterior probability.
- "The clusters we get are not the classes we want".

Unsupervised classification/clustering

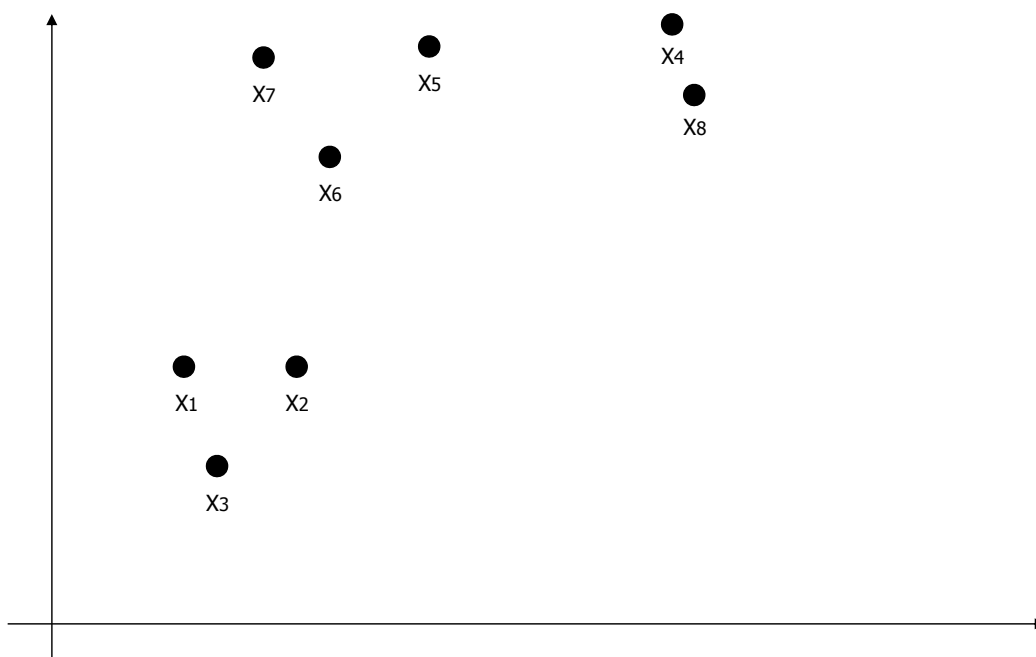
- Divide the data into clusters based on similarity (or dissimilarity)
- Similarity or dissimilarity is based on distance measures (sometimes called proximity measures)
 - Euclidean distance, Mahalanobis distance etc.
- Two main approaches to clustering
 - hierarchical
 - divisive
 - agglomerative
 - non-hierarchical (sequential)
- Non-hierarchical methods are often used in image analysis

K-means clustering

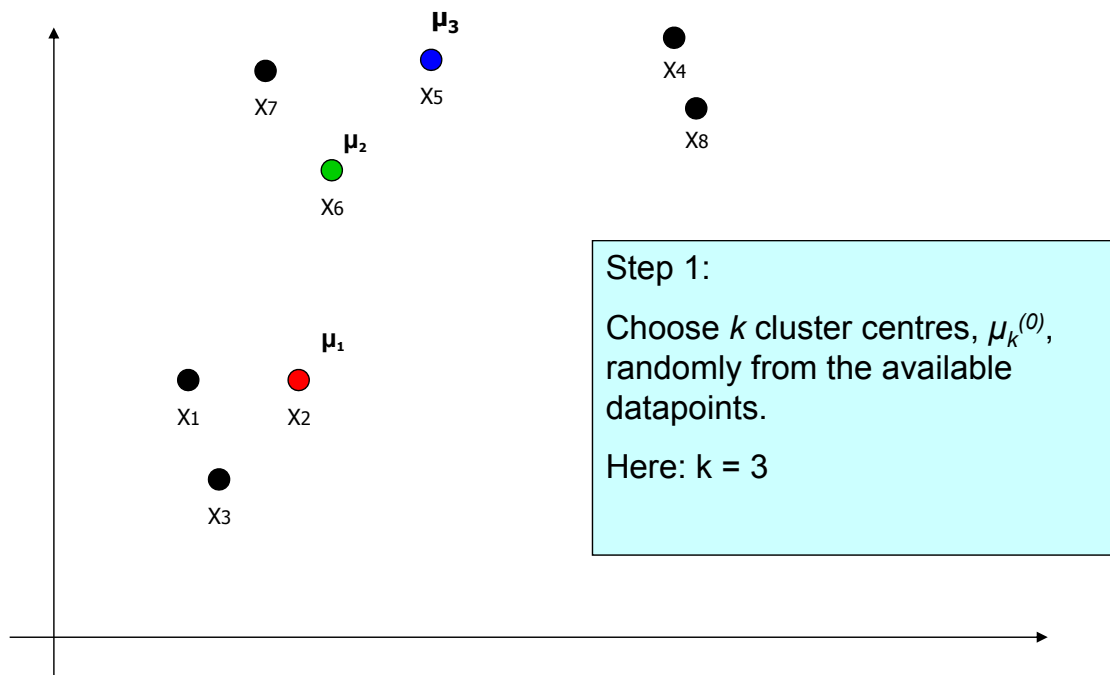
- Note: K-means algorithm normally means ISODATA, but different definitions are found in different books
- K is assumed to be known
- 1. Start with assigning K cluster centers
 - k random data points, or the first K points, or K equally spaces points
 - For $k=1:K$, Set μ_k equal to the feature vector x_k for these points.
- 2. Assign each object/pixel x_j in the image to the closest cluster center using Euclidean distance.
 - Compute for each sample the distance r^2 to each cluster center:
$$r^2 = (x_i - \mu_k)^T (x_i - \mu_k) = \|x_i - \mu_k\|^2$$
 - Assign x_j to the closest cluster (with minimum r value)
- 3. Recompute the cluster centers based on the new labels.
- 4. Repeat from 2 until #changes < limit.

ISODATA K-means: splitting and merging of clusters are included in the algorithm

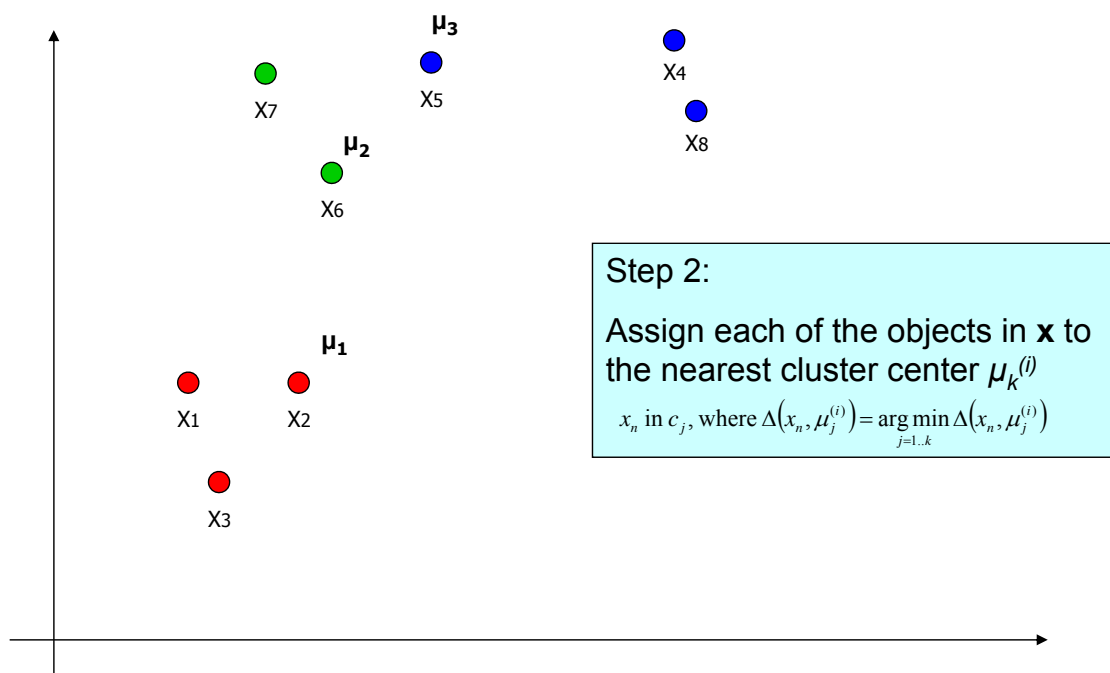
k-means example



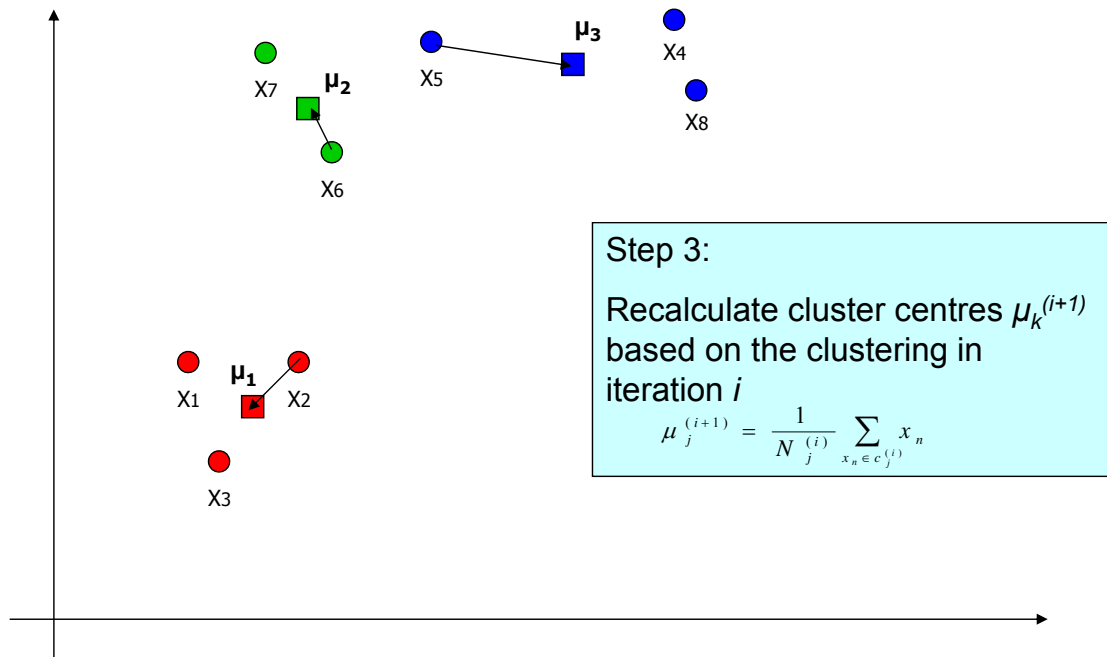
k-means example



k-means example



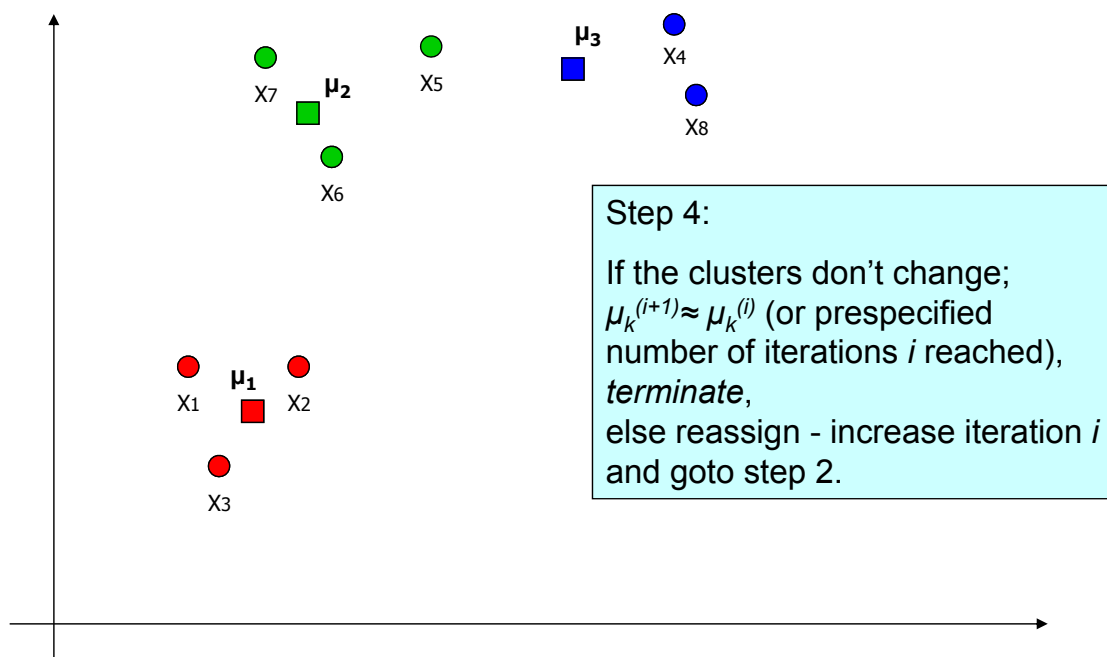
k-means example



INF 4300

35

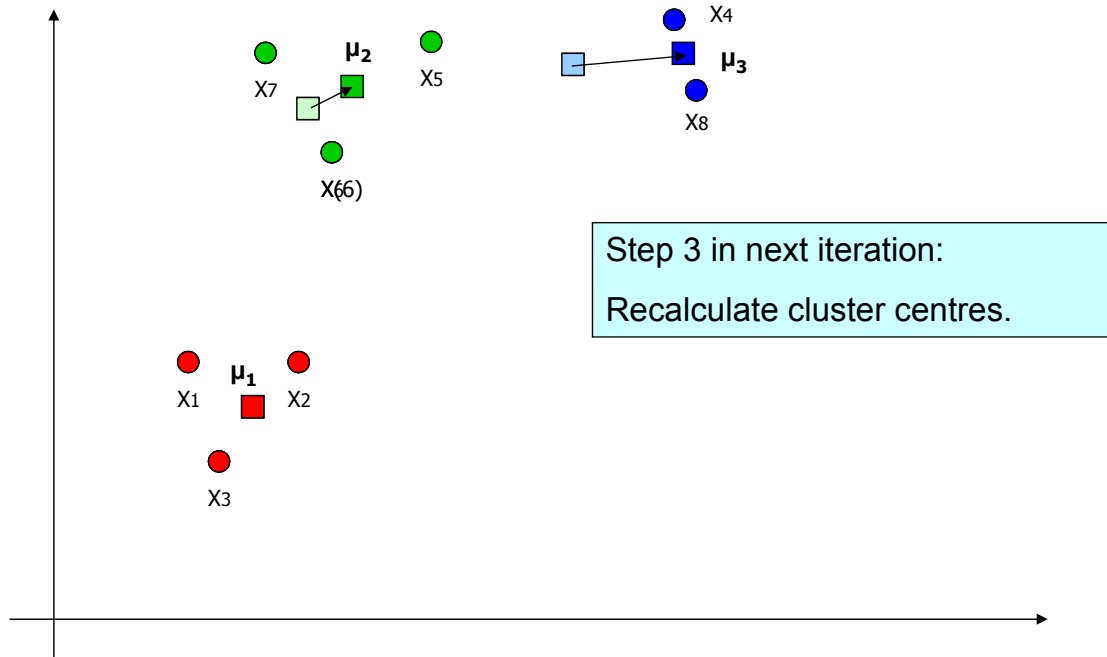
k-means example



INF 4300

36

k-means example



INF 4300

37

k-means variations

- The generic algorithm has many improvements
 - ISODATA – allow for merging and splitting of clusters
 - Among other things, this seeks to improve an initial “bad” choice of k
 - k-medians is another variation
 - k-means optimizes a probabilistic model

INF 4300

38

How do we determine k?

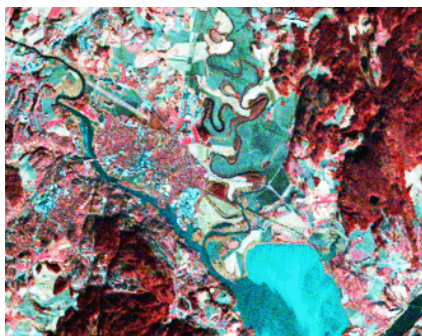
- The number of natural clusters in the data rarely corresponds to the number of information classes of interest.
- Cluster validity indices can give indications of how many clusters there are.
- Use cluster merging or splitting tailored to the application.
- Rule of thumb for practical image clustering:
 - start with approximately twice as many clusters as expected information classes
 - determine which clusters correspond to the information classes
 - split and merge clusters to improve.

INF 4300

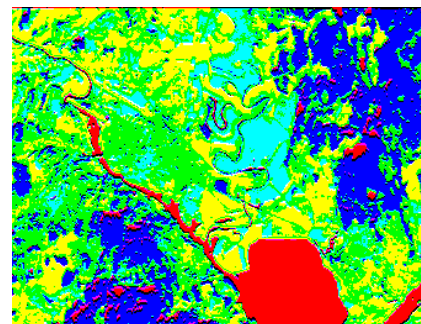
39

Example: K-means clustering

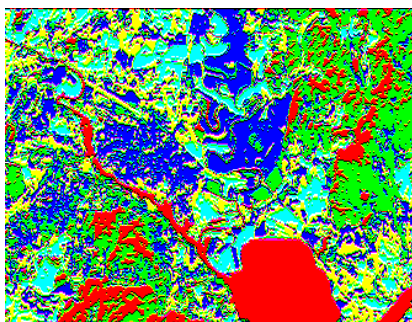
Original



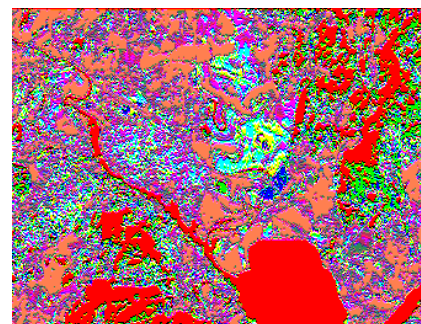
Supervised
4 classes



Kmeans
K=5



Kmeans
K=10



INF 4300

40

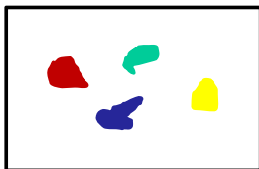
Learning goals for this lecture

- Understand how different measures of classification accuracy work:
 - Confusion matrix
 - Sensitivity/specificity/TP/TN/FP/FN
 - Average classification accuracy
- Be familiar with the curse of dimensionality and the importance of selecting few, but good features
- Know simple forward and backward feature selection.
- Understand kNN-classification
- Understand the difference between supervised and unsupervised classification
- Understand the Kmeans-algorithm.

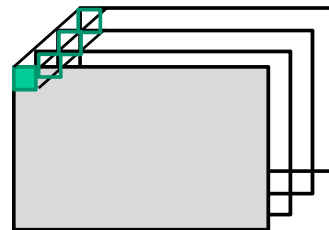
INF 4300

41

-
- \mathbf{x}_i – feature vector for pixel i
 - ω_i – The class label for pixel i
 - K – the number of classes given in the training data



Mask with training pixels



Multiband image with n spectral channels or features

$$p(\mathbf{x} | \omega_s) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}_s|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_s)^t \boldsymbol{\Sigma}_s^{-1} (\mathbf{x} - \boldsymbol{\mu}_s) \right]$$

INF 4300

42