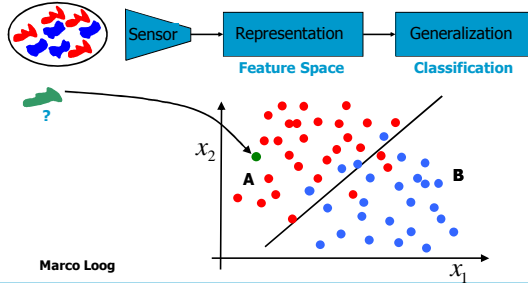


Classifiers



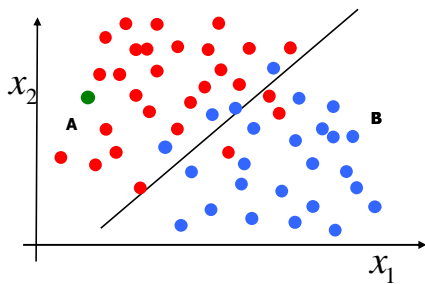
Marco Loog

ASCI A1 : Advanced Pattern Recognition

The Problem of Classification

- Learn a decision rule, from finite training data, that assigns a new object x to one of K classes
- Rule aims to minimize the error of classification
- Typically there is no perfect, i.e., zero-error, rule
- Decision rule is obtained by training a classifier

Ill-posed Problem



Some Principles

- Parametric classifiers vs. nonparametric classifiers
- Linear vs. nonlinear
- Generative classifiers : focus on each class separately, model class conditional densities [likelihoods] and reason about discrimination
- Discriminative classifiers : focus on discrimination directly, model decision function [or posterior probabilities]

Normal Density-based Classifiers

$$p(\mathbf{x} | \omega_i) = \frac{1}{\sqrt{2\pi^k \det(\Sigma_i)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_i)^\top \Sigma_i^{-1} (\mathbf{x} - \mu_i)\right) = \Phi(\mathbf{x}, \mu_i, \Sigma_i)$$

Normal Density-based Classifiers

Bayes classifier $S(\mathbf{x}) = p(\mathbf{x} | A)p(A) - p(\mathbf{x} | B)p(B)$

logs don't matter $R(\mathbf{x}) = \log(p(\mathbf{x} | A)p(A)) - \log(p(\mathbf{x} | B)p(B))$

$R(\mathbf{x}) = \log(p(\mathbf{x} | A)) - \log(p(\mathbf{x} | B)) + \log(p(A) / p(B))$

normal distribution $p(\mathbf{x} | A) = \frac{1}{\sqrt{2\pi^k \det(\Sigma_A)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_A)^\top \Sigma_A^{-1} (\mathbf{x} - \mu_A)\right)$

$\log(p(\mathbf{x} | A)) = -\frac{1}{2}(\mathbf{x} - \mu_A)^\top \Sigma_A^{-1} (\mathbf{x} - \mu_A) - \log(\sqrt{2\pi^k \det(\Sigma_A)})$

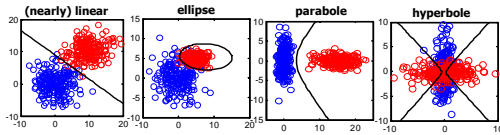
substitute $R(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_A)^\top \Sigma_A^{-1} (\mathbf{x} - \mu_A) + \frac{1}{2}(\mathbf{x} - \mu_B)^\top \Sigma_B^{-1} (\mathbf{x} - \mu_B) + c$

Quadratic expression

Normal Density-based Classifiers

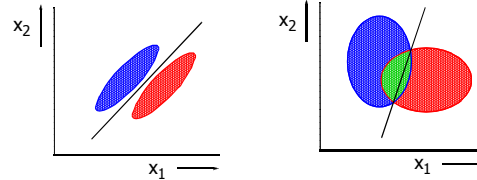
$$R(x) = -\frac{1}{2}(x - \hat{\mu}_A)^T \hat{\Sigma}_A^{-1} (x - \hat{\mu}_A) + \frac{1}{2}(x - \hat{\mu}_B)^T \hat{\Sigma}_B^{-1} (x - \hat{\mu}_B) + c$$

$$c = \log(\rho(A) / \rho(B)) + \frac{1}{2} \log[\det(\hat{\Sigma}_B) / \det(\hat{\Sigma}_A)]$$



TU Delft

Normal Density-based Classifiers



Normal distributions with equal covariance matrices Σ are optimally separated by a linear classifier

$$S(x) = (\mu_A - \mu_B)^T \Sigma^{-1} x + c$$

Optimal classifier for normal distributions with unequal covariance matrices Σ_A and Σ_B can be approximated by:

$$S(x) = (\mu_A - \mu_B)^T (\rho(A)\Sigma_A + \rho(B)\Sigma_B)^{-1} x + c$$

TU Delft

Linear Discriminants

- General form : $S(x) = w^T x + w_0$
- E.g. ldc / LDA : normal density-based linear classifier

$$S(x) = (\mu_A - \mu_B)^T (\rho(A)\Sigma_A + \rho(B)\Sigma_B)^{-1} x + c$$

- E.g. fisherc : minimum variance linear classifier

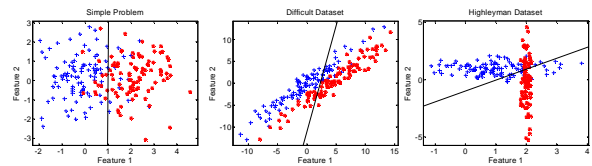
$$S(x) = (\mu_A - \mu_B)^T \Sigma^{-1} x + c$$

- Solves linear regression to indicator variables

TU Delft

Nearest Mean Classifier

$$S(x) = (\mu_A - \mu_B)^T x - (\mu_A - \mu_B)^T (\mu_A + \mu_B) / 2$$



TU Delft

Parzen Classifier

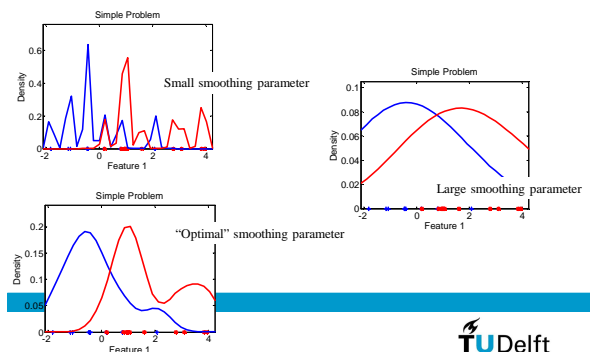
- Substitute Parzen density estimates

$$p(x | \omega) = \frac{1}{n} \sum_{x_i \in \omega} \Phi\left(\frac{x_i - x}{h}\right)$$

- Parzenc : optimize h for classification
- Parzencd : optimize h for density estimation per class

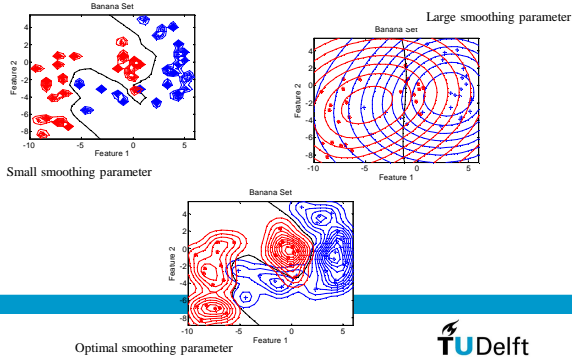
TU Delft

Parzen Classifier : 1D Example



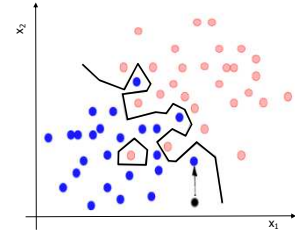
TU Delft

Parzen Classifier : 2D Example



TU Delft

Nearest Neighbor Rule

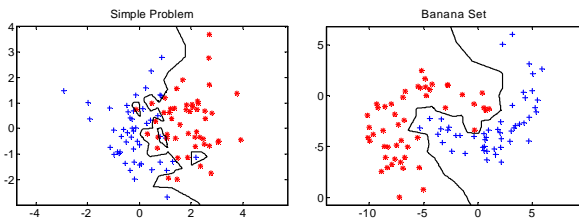


- The scaling issue...

TU Delft

More Nearest Neighbor

- Good for [almost] separable classes
- Useful for non-linear decision functions
- No training time, but long execution time
- All data should be stored



Asymptotics

- Nearest neighbor rule will, asymptotically, not perform worse than twice the best possible classifier

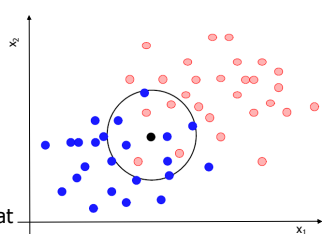
$$\epsilon^* \leq \epsilon_{1NN} \leq 2\epsilon^*(1 - \epsilon^*) \leq 2\epsilon^*$$

- Indeed, 1NN often performs rather well and provides good baseline results

TU Delft

k-Nearest Neighbor Classifier

- More smooth
- Less local

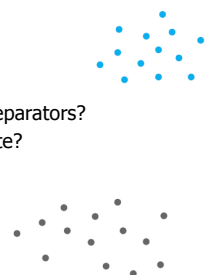


- Generally, one can show that

$$\epsilon^* \leq \epsilon_{kNN} \leq (1 + \sqrt{2/k})\epsilon^*$$

TU Delft

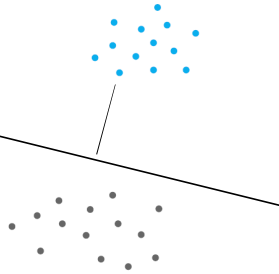
- How many perfect linear separators?
- What would be your favorite?



TU Delft

Support Vector Machines

- What's a margin?



- And how does one maximize it?



Naïve Bayes

- Assume features independent given class label
- Estimate class-conditional densities by
- Bayes rule does the rest...

$$p(x|A) = \prod_{i=1}^k p(x_i | A)$$

$$p(x|B) = \prod_{i=1}^k p(x_i | B)$$



Logistic Regression / Classification

- Based on assumption that logarithm of posterior odds is linear

$$\log\left(\frac{p(A|x)}{p(B|x)}\right) = \log\left(\frac{p(A)p(x|A)}{p(B)p(x|B)}\right) = w^T x + w_0$$

- One can derive that

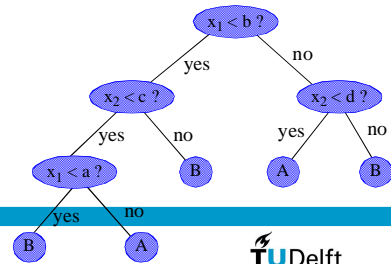
$$p(A|x) = \frac{p(A)p(x|A)}{p(A)p(x|A) + p(B)p(x|B)} = \frac{e^{w^T x + w_0}}{1 + e^{w^T x + w_0}} = \frac{1}{1 + e^{-w^T x - w_0}}$$

$$p(B|x) = \frac{p(B)p(x|B)}{p(A)p(x|A) + p(B)p(x|B)} = \frac{1}{1 + e^{w^T x + w_0}}$$

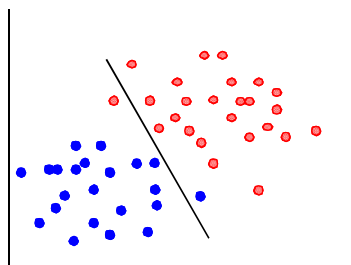
- Product over all terms in training set is optimized



Decision Trees



Perceptron



- Linear classifier
- Trained through simple update rule

$$w^{k+1} = w^k + \alpha \Delta w(w^k, x_i) = w^k + \alpha x_i$$



Somewhat more General View

- Given class M of functions / models m parameterized by θ Model $m(\cdot|\theta)$ takes x as input and provides posterior, label, whatever as output
- Given training data $(x_s \text{ and } y_s)$, find that θ that minimizes the empirical risk :

$$\sum_i \ell(m(x_i|\theta), y_i)$$

- Loss ℓ defines how to penalize deviations of model prediction for x from true value y



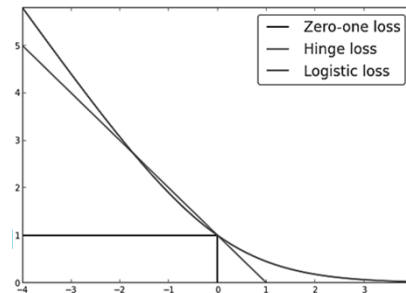
Some Losses

- 0/1 : $\begin{cases} 0, m(x|\theta) = y \\ 1, \text{otherwise} \end{cases}$
 - Squared : $(m(x|\theta) - y)^2$
 - Logistic : $\log(1 + \exp(-y m(x|\theta)))$
 - Soft-margin / hinge : $\max(0, 1 - y m(x|\theta))$
- What loss is often used in generative modeling?



Why Approximate?

- In fact, we would like to optimize 0/1 in classification



Additional Remarks

- Rather than empirical risk, one often optimizes regularized / penalized / ... risk

$$\sum_i \ell(m(x_i|\theta), y_i) + R(\theta)$$

- A priori, R -term biases solution to simple ones
- Needed when models too flexible / little training data
- Not all classifiers can be formulated in terms of above formalism... Or?



PRTools : Some Generative Classifiers

- NMC - nearest mean classifier
- NMSC - nearest mean scaled classifier
- FISHERC - Fisher linear discriminant
- LDC - linear discriminant
- QDC - quadratic discriminant
- UDC - quadratic discriminant diagonal covariance
- MOGC - mixture of Gaussians classification
- PARZENC - Parzen classifier
- NAIVEBC - naïve Bayes classifier



PRTools : Some Discriminative Classifiers

- TREEC - binary decision tree classifier
- BPXNC - backpropagation feed forward neural network
- LMNC - Levenberg-Marquardt feed forward neural net
- PERLC - linear perceptron
- RBNC - radial basis neural network classifier
- SUBSC - subspace classifier
- SVC - support vector machine
- KNNC - k-nearest neighbor rule

