



OPPIA. European Social Fund
Prague & EU: We invests in your future.

LEARNING & LINEAR CLASSIFIERS

J. Matas

Czech Technical University, Faculty of Electrical Engineering
Department of Cybernetics, Center for Machine Perception
121 35 Praha 2, Karlovo nám. 13, Czech Republic

matas@fel.cvut.cz, <http://cmp.felk.cvut.cz>

LECTURE PLAN

- ◆ The problem of classifier design.
- ◆ Learning in pattern recognition.
- ◆ Linear classifiers.
- ◆ Perceptron algorithms.
- ◆ Optimal separating plane with the Kozinec algorithm.

The object of interest is characterised by observable properties $x \in X$ and its class membership (unobservable, hidden state) $k \in K$, where X is the space of observations and K the set of hidden states.

The objective of classifier design is to find a function $q^* : X \rightarrow K$ that has some optimal properties.

Bayesian decision theory solves the problem of minimisation of risk

$$R(q) = \sum_{x,k} W(q(x), k) p(x, k)$$

given the following quantities:

- ◆ $p(x, k), \forall x \in X, k \in K$ – the statistical model of the dependence of the observable properties (measurements) on class membership
- ◆ $W(q(x), k)$ the loss of decision $q(x)$ if the true class is k

Do you know the solution for the 0-1 loss function?

Non-Bayesian decision theory solves the problem if $p(x|k), \forall x \in X, k \in K$ are known, but $p(k)$ are unknown (or do not exist). Constraints or preferences for different errors depend on the problem formulation.

Do you know any non-bayesian problem formulations?

However, in applications typically:

- ◆ none of the probabilities are known! The designer is only given a **training multiset** $T = \{(x_1, k_1) \dots (x_L, k_L)\}$, where L is the length (size) of the training multiset.
- ◆ the desired properties of the classifier $q(x)$ are known

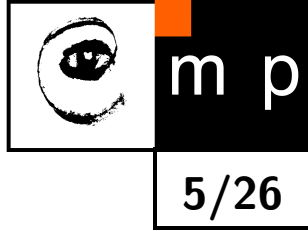
How would you proceed ?

CLASSIFIER DESIGN via PARAMETER ESTIMATION



- ◆ Assume $p(x, k)$ have a particular form, e.g. Gaussian (mixture), piece-wise constant, etc., with a finite (i.e. small) number of parameters Θ_k .
 - ◆ Estimate the parameters from the using training set T
 - ◆ Solve the classifier design problem (e.g. risk minimisation), substituting the estimated $\hat{p}(x, k)$ for the true (and unknown) probabilities $p(x, k)$
- ? : What estimation principle should be used?
- What estimation paradigms do you know?*
- : There is no direct relationship between known properties of estimated $\hat{p}(x, k)$ and the properties (typically the risk) of the obtained classifier $q'(x)$
 - : If the true $p(x, k)$ is not of the assumed form, $q'(x)$ may be arbitrarily bad, even if the size of training set L approaches infinity!
 - + : Implementation is often straightforward, especially if parameters Θ_k for each class are assumed independent.
 - + : Performance on training data can be predicted by crossvalidation.

LEARNING in STATISTICAL PATTERN RECOGNITION



- ◆ Choose a class Q of decision functions (classifiers) $q : X \rightarrow K$.
- ◆ Find $q^* \in Q$ minimising some criterion function on the training set that approximates the risk $R(q)$ (which cannot be computed).
- ◆ Learning paradigm is defined by the criterion function:

Empirical risk (training set error) minimization. True risk approximated

$$R_{emp}(q_{\Theta}(x)) = \frac{1}{L} \sum_{i=1}^L W(q_{\Theta}(x_i), k_i),$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} R_{emp}(q_{\Theta}(x))$$

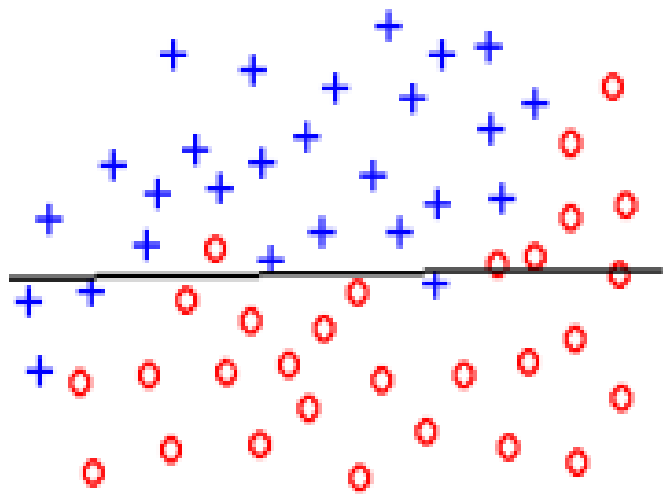
Examples: Perceptron, Neural nets (Back-propagation), etc.

Structural risk minimization.

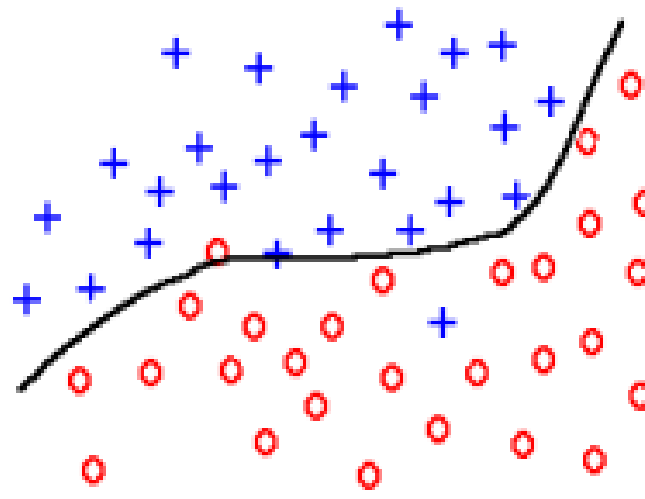
Example: SVM (Support Vector Machines).

OVERFITTING AND UNDERFITTING

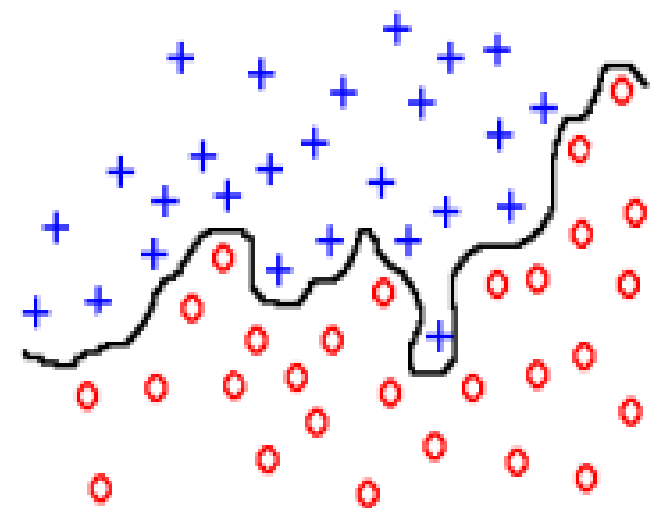
- ◆ How rich class \mathcal{Q} of classifiers $q_{\Theta}(x)$ should be used?
- ◆ The problem of generalization is a key problem of pattern recognition: a small empirical risk R_{emp} need not imply a small true expected risk R !



underfit



fit



overfit

STRUCTURAL RISK MINIMIZATION PRINCIPLE (1)



We would like to minimise the risk

$$R(q) = \sum_{x,k} W(q_{\Theta}(x), k) p(x, k)$$

but $p(x, k)$ is unknown.

Vapnik and Chervonenkis proved a remarkable inequality

$$R(q) \leq R_{emp}(q) + R_{str} \left(h, \frac{1}{L} \right) ,$$

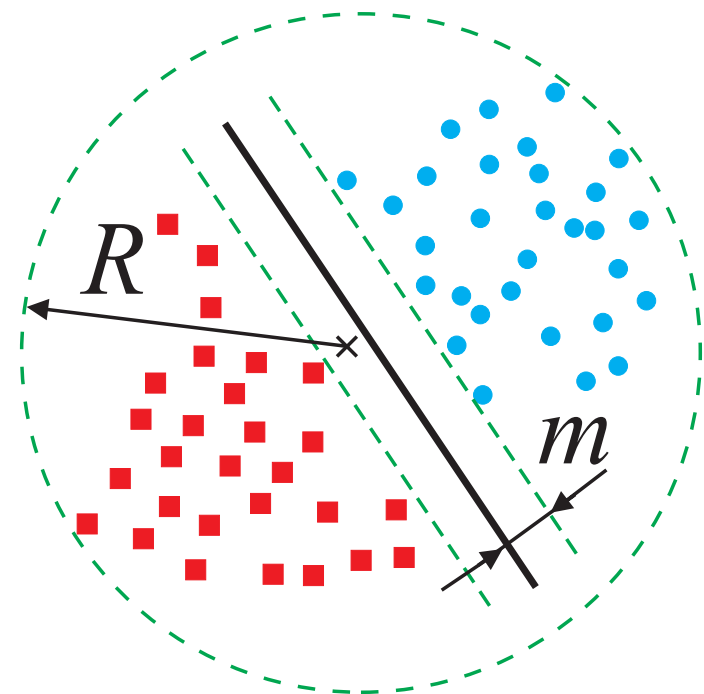
where h is VC dimension (capacity) of the class of strategies Q .

Notes:

- + R_{str} does not depend on the unknown $p(x, k)$!!
- + R_{str} known for some classes of Q , e.g. linear classifiers.

STRUCTURAL RISK MINIMIZATION PRINCIPLE (2)

- ◆ There are more types of upper bounds on R .
E.g. for linear discriminant functions



VC dimension (capacity)

$$h \leq \frac{R^2}{m^2} + 1$$

- ◆ Examples of learning algorithms: SVM or ε -Kozinec.

$$(w^*, b^*) = \operatorname{argmax}_{w, b} \min \left(\min_{x \in X_1} \frac{\langle w, x \rangle + b}{|w|}, \min_{x \in X_2} \frac{\langle w, x \rangle + b}{|w|} \right).$$

Is then empirical risk minimisation = minimisation of training set error, e.g. neural networks with backpropagation, dead ? No!

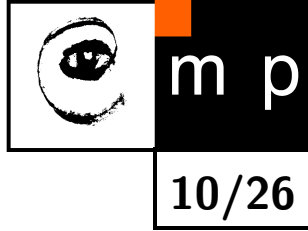
– R_{str} may be so large that the upper bound is useless.

Find a tighter bound and you will be famous! It is not impossible!

- + Vapnik's theory justifies using empirical risk minimisation on classes of functions with VC dimension.
- + Vapnik suggests learning with progressively more complex classes Q .
- + Empirical risk minimisation is computationally hard (impossible for large L). Most classes of decision functions Q where empirical risk minimisation (at least local) can be efficiently organised are often useful.

Where does the nearest neighbour classifier fit in the picture?

WHY ARE LINEAR CLASSIFIERS IMPORTANT?



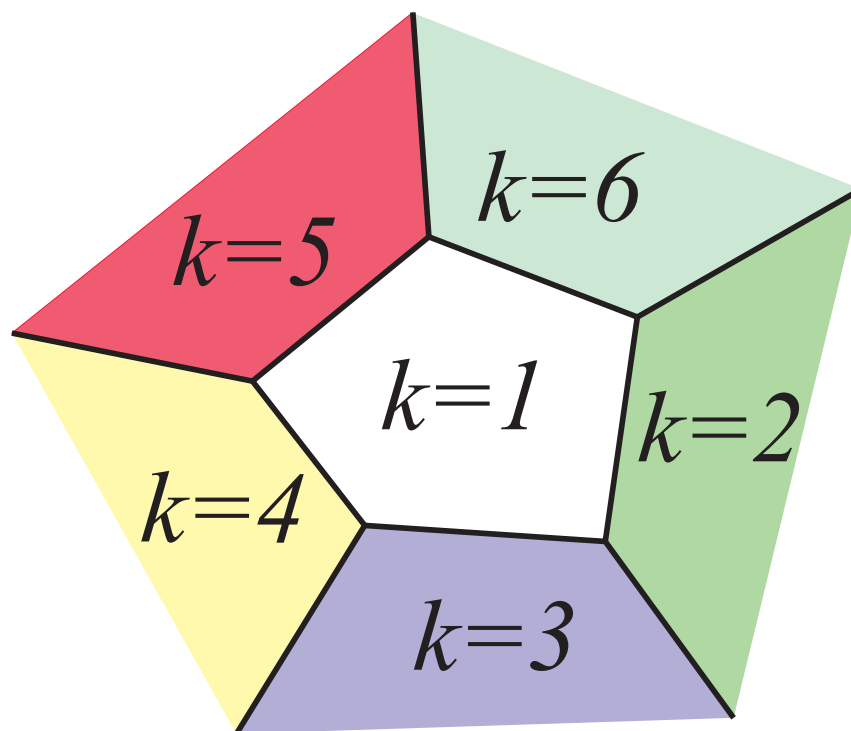
- ◆ For some statistical models, the Bayesian or non-Bayesian strategy is implemented by a linear discriminant function.

You should know an example!?

- ◆ Capacity (VC dimension) of linear strategies in an n -dimensional space is $n + 2$. Thus, the learning task is well-posed, i.e., strategy tuned on a finite training multiset does not differ much from correct strategy found for a statistical model.
- ◆ There are efficient learning algorithms for linear classifiers.
- ◆ Some non-linear discriminant functions can be implemented as linear after the feature space transformation.

LINEAR DISCRIMINANT FUNCTION $q(x)$

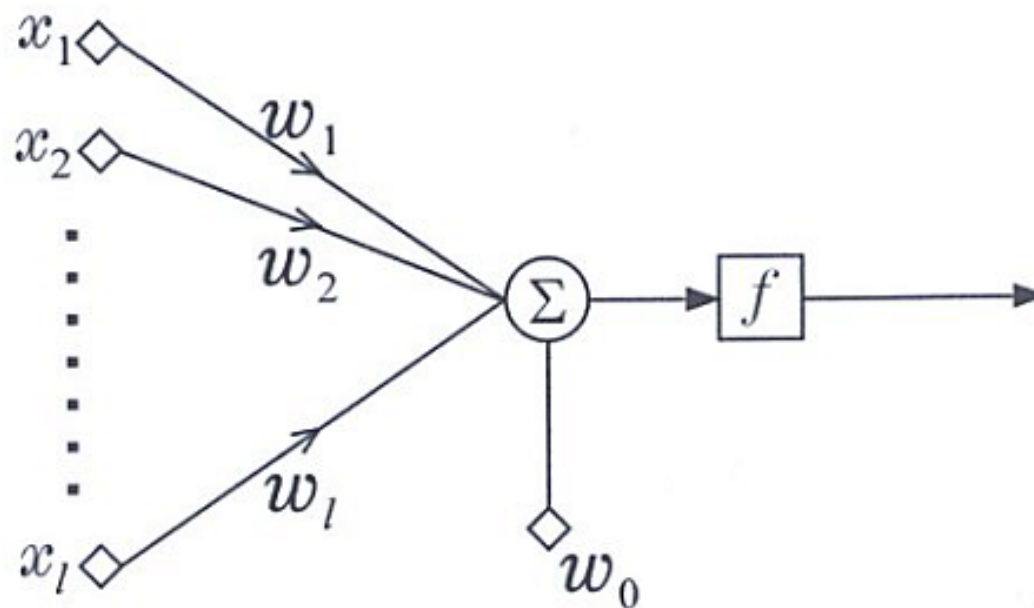
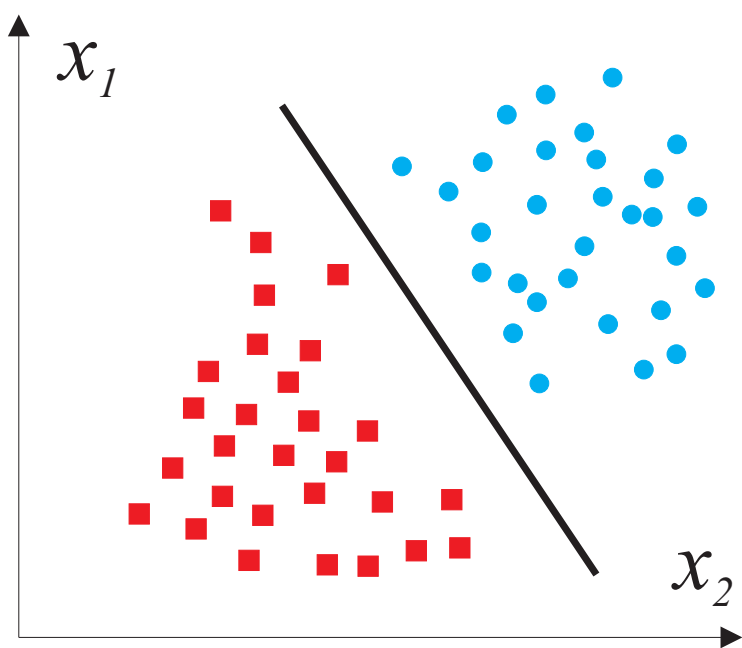
- ◆ $f_j(x) = \langle w_j, x \rangle + b_j$, where $\langle \rangle$ denotes a scalar product.
- ◆ A strategy $j = \operatorname{argmax}_j f_j(x)$ divides X into $|K|$ convex regions.



DICHOTOMY, TWO CLASSES ONLY

$|K| = 2$, i.e. two hidden states (typically also classes)

$$q(x) = \begin{cases} k = 1, & \text{if } \langle w, x \rangle + b \geq 0, \\ k = -1, & \text{if } \langle w, x \rangle + b < 0. \end{cases}$$



PERCEPTRON LEARNING: Formulation

Input: $T = \{(x_1, k_1) \dots (x_L, k_L)\}, k \in \{-1, 1\}$

Output: a weight vector w , offset b , satisfying, $\forall j \in \{1..L\}$:

$$\langle w, x_j \rangle + b \geq 0 \quad \text{if } k_j = 1,$$

$$\langle w, x_j \rangle + b < 0 \quad \text{if } k_j = -1$$

equivalently, multiplying both inequalities by k_j ,

$$\langle w, k_j x_j \rangle + k_j b \geq 0$$

or even simpler

$$\langle w', k_j x'_j \rangle \geq 0,$$

where $x' = [x \quad 1], w' = [w \quad b]$

PERCEPTRON LEARNING: Simplified Formulation

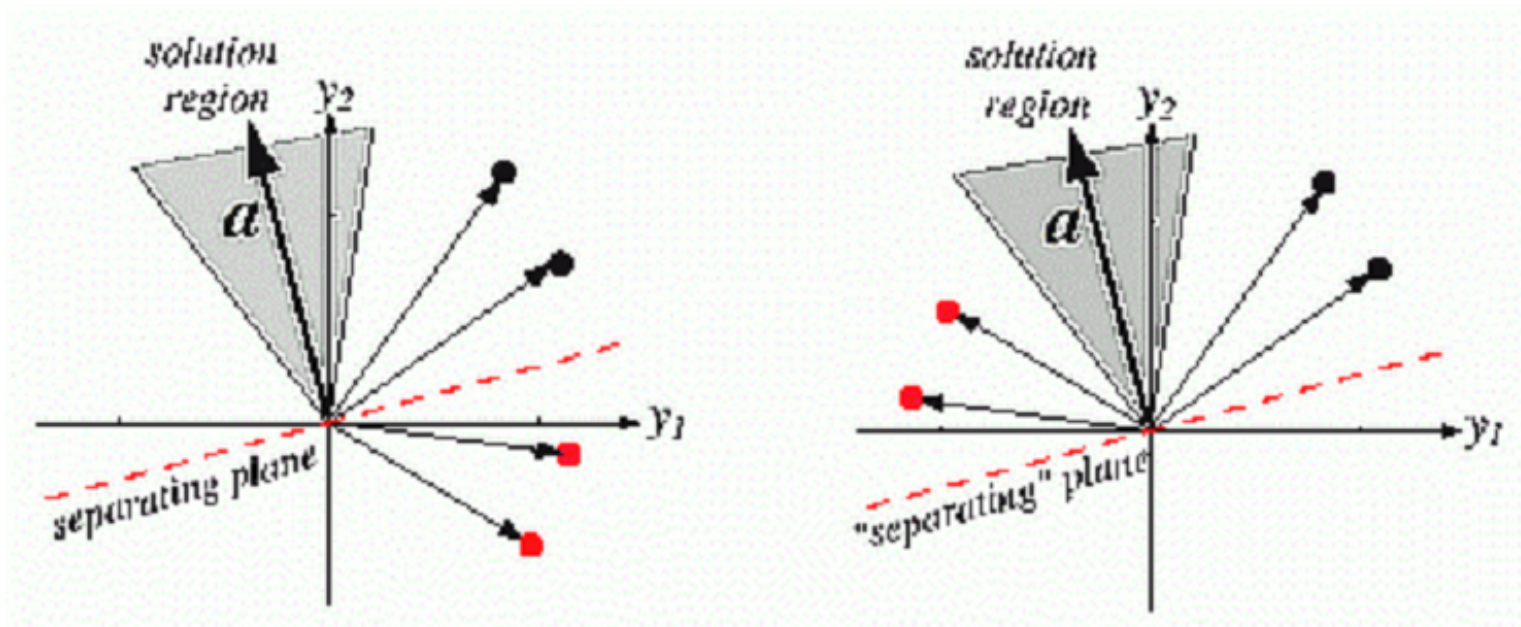
To simplify notation, we reformulate the problem.

Input: $T = \{x''_1, \dots, x''_L\}$ where $x''_j = k_j [x_j \ 1]$

Output: a weight vector $w' = [w \ b]$, such that :

$$\langle w', x''_j \rangle \geq 0, \forall j \in \{1..L\}$$

We drop the primes and go back to w, x notation. The vector w has the offset b as last element, x has an extra 1 and has been multiplied by k .



PERCEPTRON LEARNING: THE ALGORITHM

Input: $T = \{x_1, \dots, x_L\}$

Output: a weight vector w

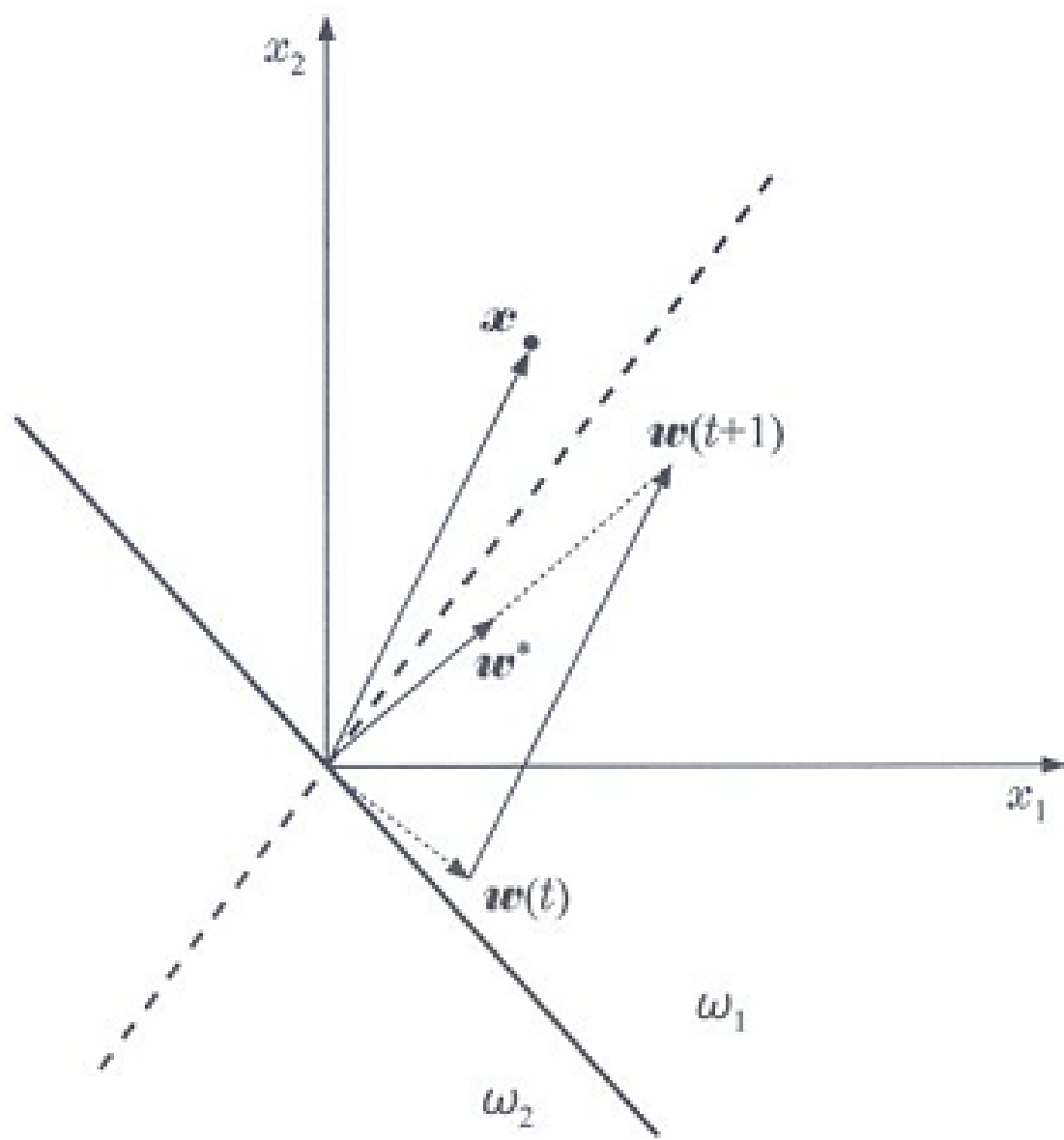
Perceptron algorithm, (Rosenblat 1962):

1. $w_1 = 0$.
2. A wrongly classified observation x_j is sought, i.e.,
 $\langle w_t, x_j \rangle < 0, j \in \{1..L\}$.
3. If there is no misclassified observation then the algorithm terminates otherwise

$$w_{t+1} = w_t + x_j .$$

4. Goto 2.
-

PERCEPTRON: WEIGHT UPDATE



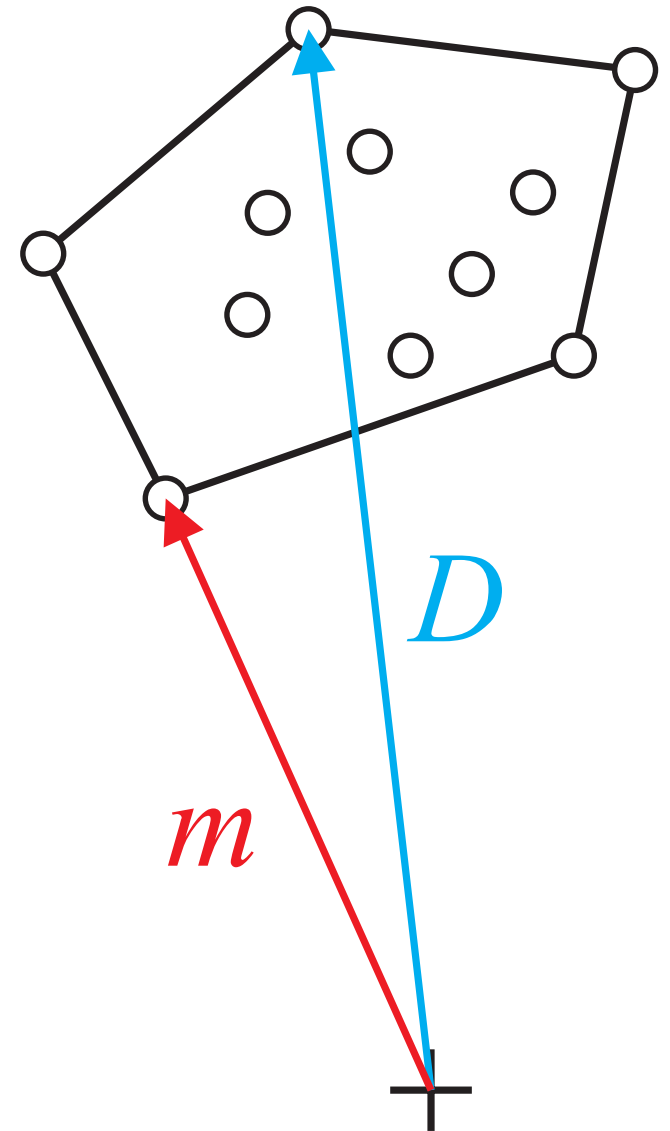
NOVIKOFF THEOREM

If the data are linearly separable then there exists a number $t^* \leq \frac{D^2}{m^2}$, such that the vector w_{t^*} satisfies the inequality

$$\langle w_{t^*}, x^j \rangle > 0, \forall j \in \{1..L\}.$$

? What if the data is not separable?

? How to terminate perceptron learning?



Perceptron algorithm, batch version, handling non-separability:

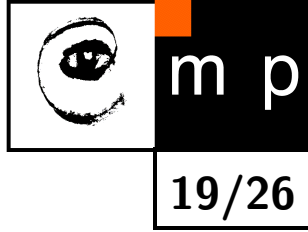
Input: $T = \{x_1, \dots, x_L\}$

Output: a weight vector w^*

1. $w_1 = 0$, $E = |T| = L$, $w^* = 0$.
2. Find all mis-classified observations $X^- = \{x \in X : \langle w_t, x \rangle < 0\}$.
3. if $|X^-| < E$ then $E = |X^-|$; $w^* = w_t$
4. if $tc(w^*, t, t_{lu})$ then terminate else $w_{t+1} = w_t + \eta_t \sum_{x \in X^-} x$
5. Goto 2.

-
- ◆ The algorithm converges with probability 1 to the optimal solution.
 - ◆ Convergence rate not known (to me).
 - ◆ Termination condition $tc(\cdot)$ is a complex function of the quality of the best solution, time since last update $t - t_{lu}$ and requirements on the solution.

PERCEPTRON LEARNING as an Optimisation problem (1)



Perceptron algorithm, batch version, handling non-separability, another perspective:

Input: $T = \{x_1, \dots, x_L\}$

Output: a weight vector w minimising

$$J(w) = |\{x \in X : \langle w_t, x \rangle < 0\}|$$

or, equivalently

$$J(w) = \sum_{x \in X : \langle w_t, x \rangle < 0} 1$$

What would the most common optimisation method, i.e. [gradient descent](#), perform?

$$w_t = w - \eta \nabla J(w)$$

The gradient of $J(w)$ is either 0 or undefined. Gradient minimisation cannot proceed.

PERCEPTRON LEARNING as an Optimisation problem (2)

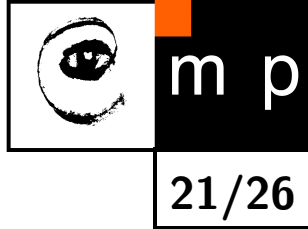
Let us redefine the cost function:

$$J_p(w) = \sum_{x \in X: \langle w, x \rangle < 0} \langle w, x \rangle$$

$$\nabla J_p(w) = \frac{\partial J}{\partial w} = \sum_{x \in X: \langle w, x \rangle < 0} x$$

- ◆ The Perceptron Algorithm is a gradient descent method for $J_p(w)$!
- ◆ Learning and empirical risk minimisation is just an instance of an [optimization problem](#).
- ◆ Either gradient minimisation (backpropagation in neural networks) or convex (quadratic) minimisation (in mathematical literature called convex programming) is used.

OPTIMAL SEPARATING PLANE and THE CLOSEST POINT TO THE CONVEX HULL



The problem of optimal separation by a hyperplane

$$(1) \quad w^* = \operatorname{argmax}_w \min_j \left\langle \frac{w}{|w|}, x_j \right\rangle$$

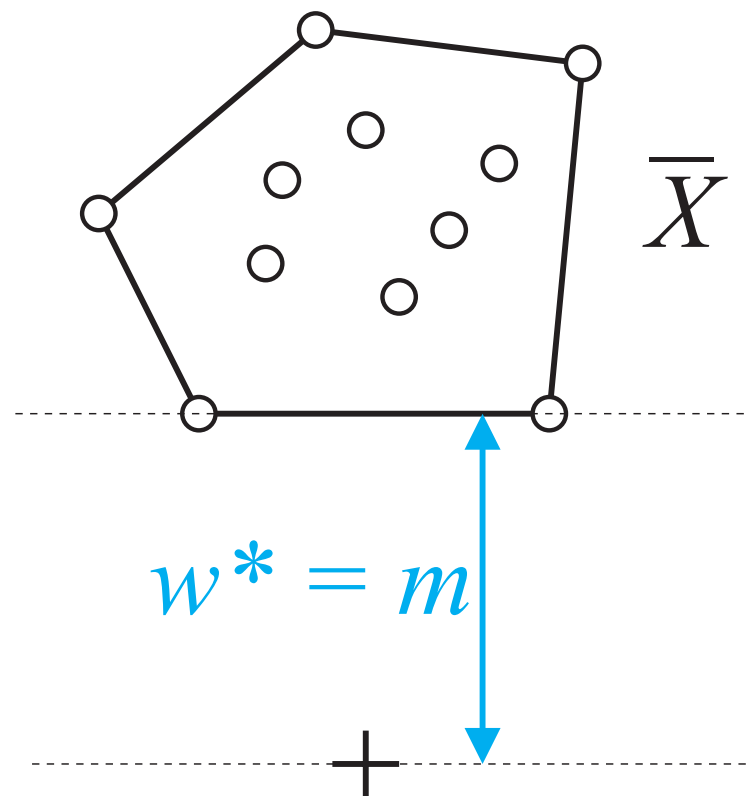
can be converted to seek for the closest point to a convex hull (denoted by the overline)

$$x^* = \operatorname{argmin}_{x \in \overline{X}} |x|$$

There holds that x^* solves also the problem (1).

Recall that the classifier that maximises separation minimises the structural risk R_{str} (page 8)!

CONVEX HULL, ILLUSTRATION



$$\min_j \left\langle \frac{w}{|w|}, x_j \right\rangle \leq m \leq |w|, w \in \bar{X}$$

lower bound

upper bound

ε -SOLUTION

- ◆ The aim is to speed up the algorithm.
- ◆ The allowed uncertainty ε is introduced.

$$|w| - \min_j \left\langle \frac{w}{|w|}, x_j \right\rangle \leq \varepsilon$$

TRAINING ALGORITHM 2 – KOZINEC (1973)

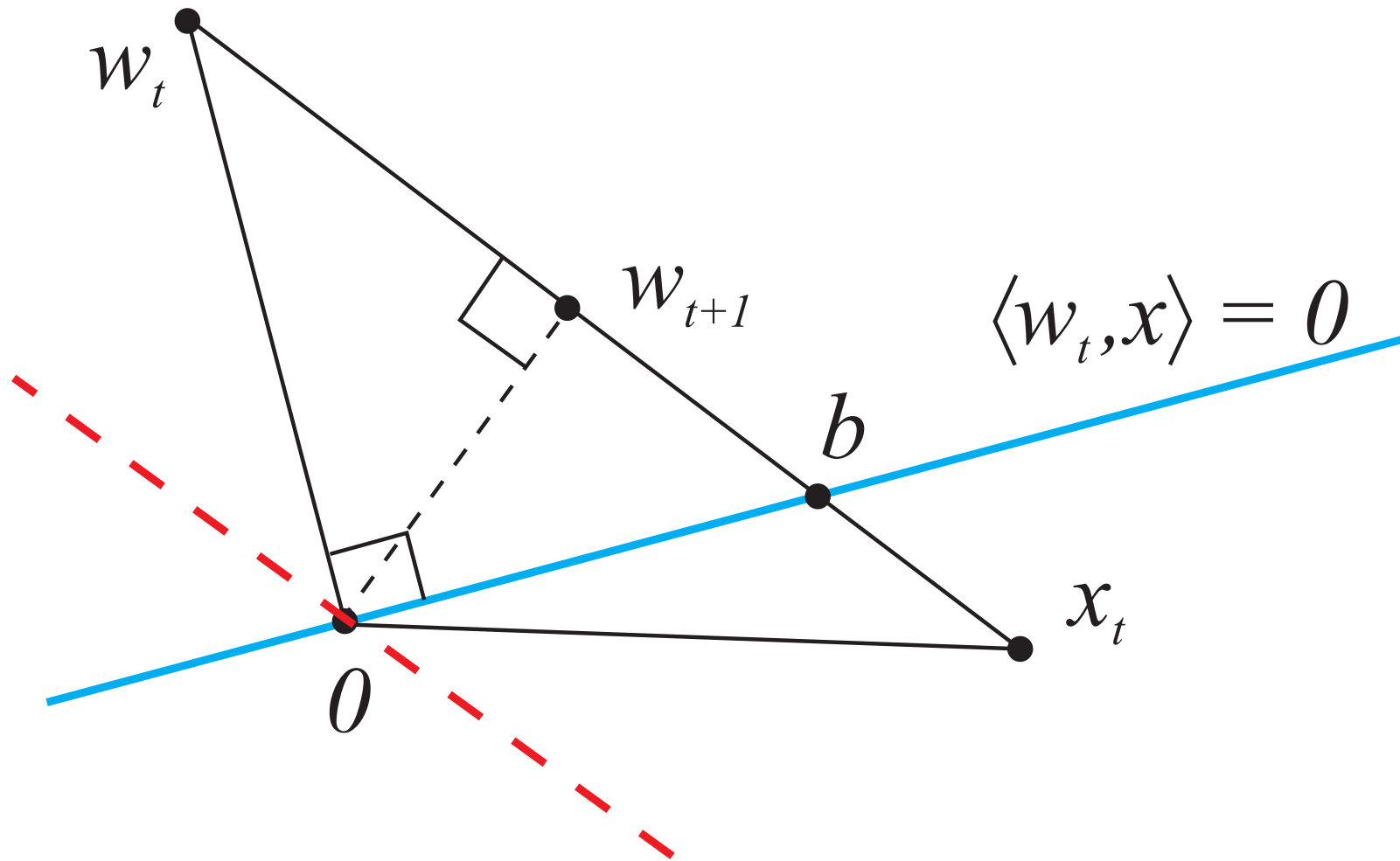
1. $w_1 = x_j$, i.e. any observation.
2. A wrongly classified observation x_t is sought, i.e.,
 $\langle w_t, x^j \rangle < b, j \in J$.
3. If there is no wrongly classified observation then the algorithm finishes otherwise

$$w_{t+1} = (1 - k) \cdot w_t + x_t \cdot k, \quad k \in \mathbb{R}.$$

where $k = \underset{k}{\operatorname{argmin}} |(1 - k) \cdot w_t + x_t \cdot k|$.

4. Goto 2.

KOZINEC, PICTORIAL ILLUSTRATION



Kozinec

KOZINEC and ϵ -SOLUTION

The second step of Kozinec algorithm is modified to:

A wrongly classified observation x_t is sought, i.e.,

$$|w^t| - \min_j \left\langle \frac{w^t}{|w^t|}, x_t \right\rangle \geq \epsilon$$

