

Computational learning and discovery



CSI 873 / MATH 689

Instructor: I. Griva

Wednesday 7:20 - 10 pm

Support vector machine

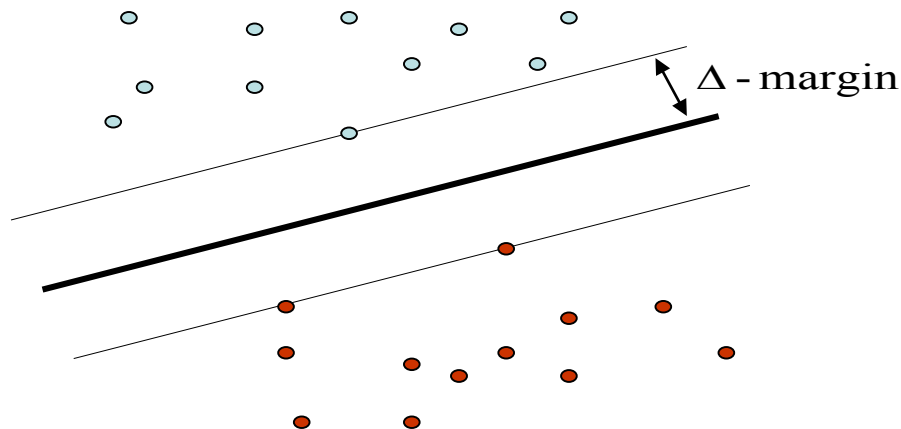
Given a set of training data

$$(x_1, y_1), \dots, (x_l, y_l), x_i \in \mathbb{R}^n, y_i \in \{+1, -1\}$$

find a function that can estimate

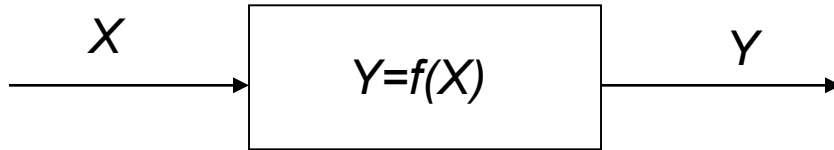
$$y_j^* \in \{+1, -1\} \text{ given new } x_j^* \in \mathbb{R}^n$$

and minimize the frequency of the future error.



Support vector machine

based on fundamentals of statistical learning theory
(Vapnik-Chervonenkis theory)



Instead of identifying the unknown function (what classical statistics does), the main goal of VC theory is to imitate the unknown function.

The key discovery of VC theory:

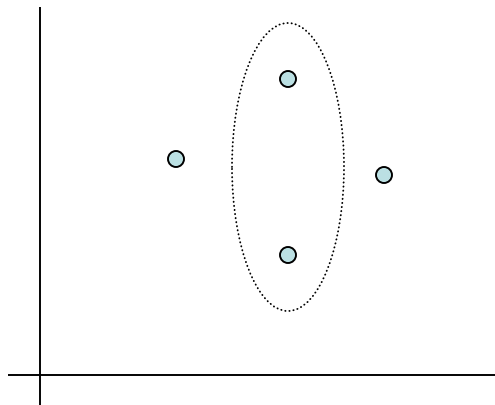
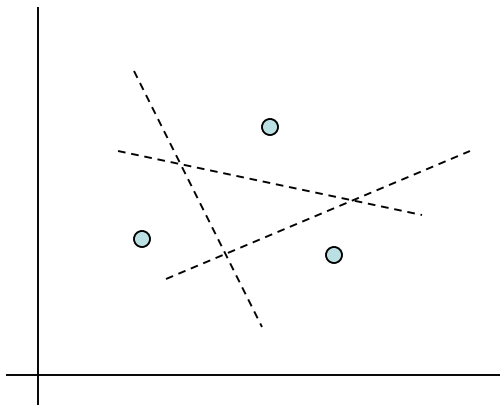
- **Two and only two factors are responsible for generalization:**
 - **One (empirical loss) defines how well the function approximates data**
 - **Another (capacity, VC dimension) defines the diversity of the set of functions from which one chooses an approximation function**
- **If VC dimension is finite, then one can achieve a good generalization. If it is not finite the generalization is impossible.**

Examples

The VC dimension of linear indicator functions

$$I(x) = \text{sgn}(x^T w + b), \quad x \in \mathbb{R}^n, w \in \mathbb{R}^n$$

is equal $n + 1$

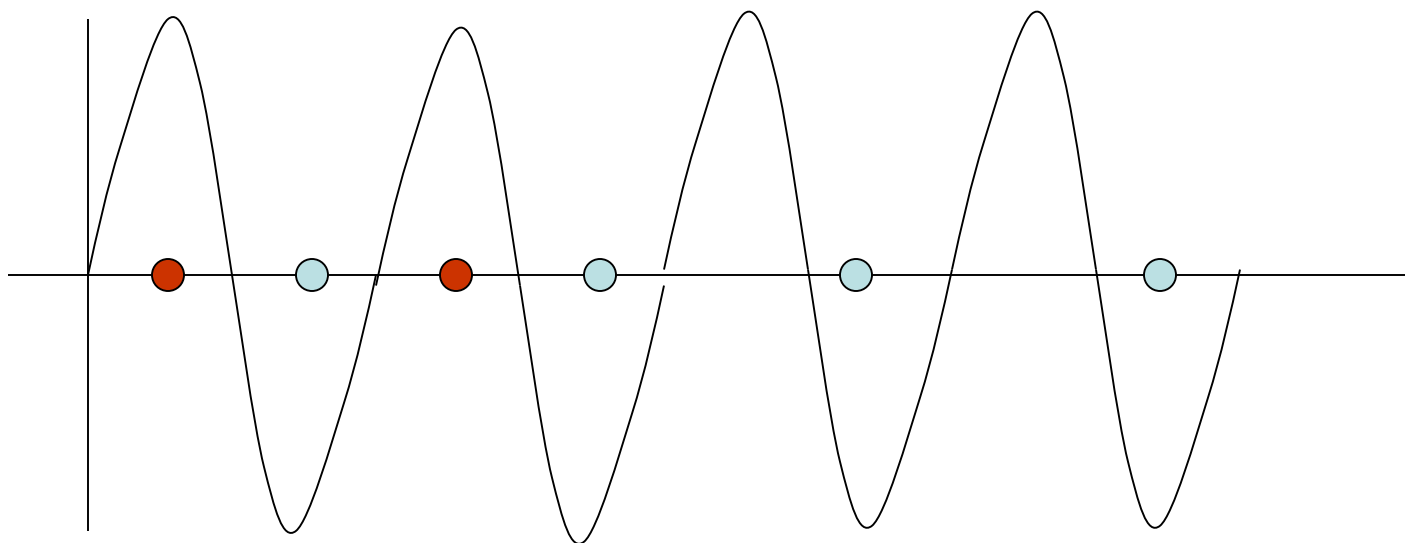


Examples

The VC dimension of the set of functions

$$I(x) = \text{sgn}(\sin ax), \quad x \in \mathbb{R}^1, w \in \mathbb{R}^1$$

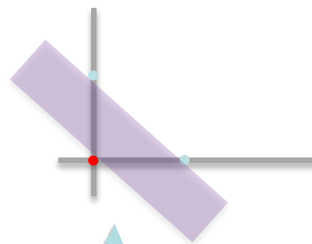
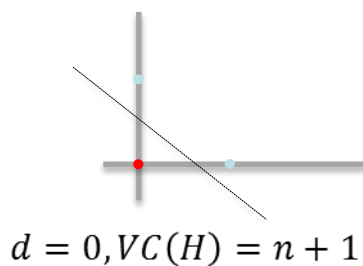
is infinity



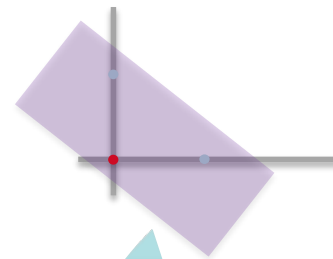
Support vector machine

Let the vector $x \in \Re^n$ belong to a sphere of radius R . Then the set of Δ - margin separating hyperplanes has a VC dimension bounded as follows

$$VC_{\text{dim}} \leq \min \left\{ \frac{R^2}{\Delta^2}, n \right\} + 1$$



This set is can be shattered by a hyperplane with a margin $d > 0$

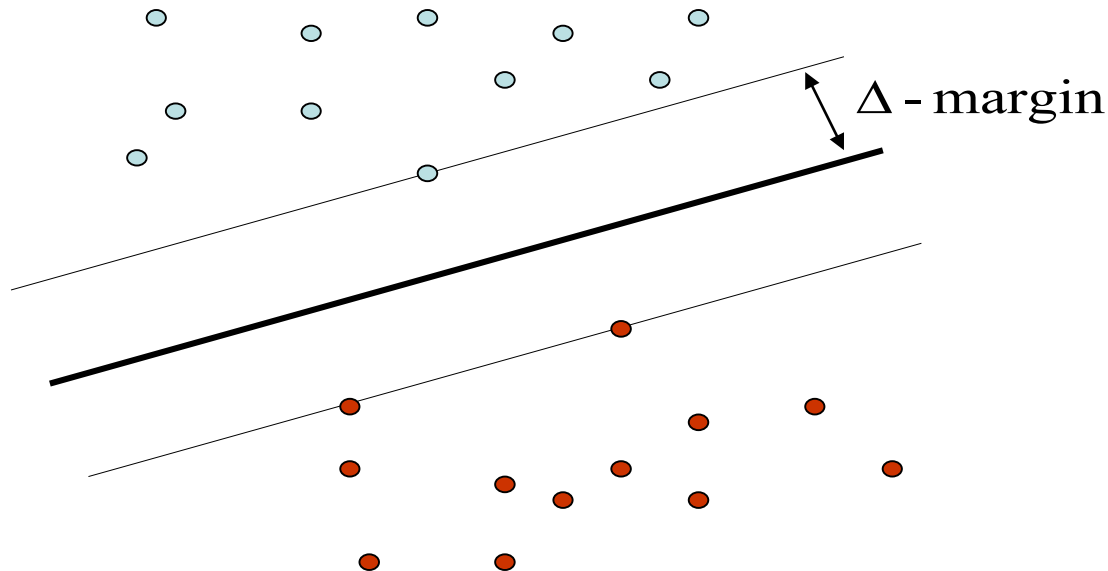


Margin d is too large to shatter this set

Support vector machine

Let the vector $x \in \Re^n$ belong to a sphere of radius R . Then the set of Δ - margin separating hyperplanes has a VC dimension bounded as follows

$$VC_{\text{dim}} \leq \min \left\{ \frac{R^2}{\Delta^2}, n \right\} + 1$$



Support vector machine

Theorem. With probability $1 - \delta$ one can assert that the probability that a test example will not be separated correctly by the Δ - margin hyperplane has the bound

$$P_{error} \leq \frac{m}{l} + \frac{\kappa}{2} \left(1 + \sqrt{1 + \frac{4m}{l\kappa}} \right),$$

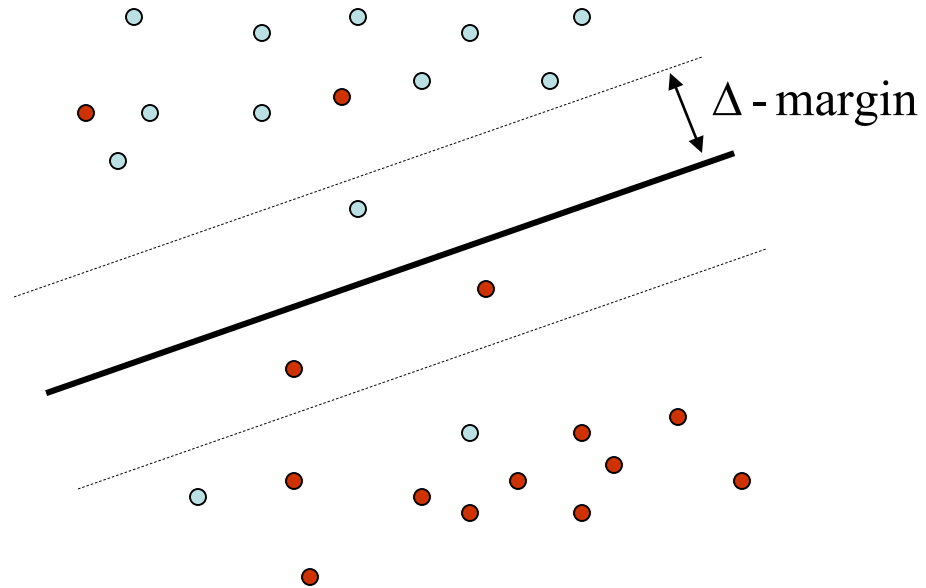
where

$$\kappa = 4 \frac{VC_{\dim} \left(\ln \frac{2l}{VC_{\dim}} + 1 \right) - \ln \frac{\delta}{4}}{l},$$

m is the number of training examples that are not separated correctly by the Δ - margin hyperplane and the VC dimension bounded as follows

$$VC_{\dim} \leq \min \left\{ \frac{R^2}{\Delta^2}, n \right\} + 1$$

Support vector machine



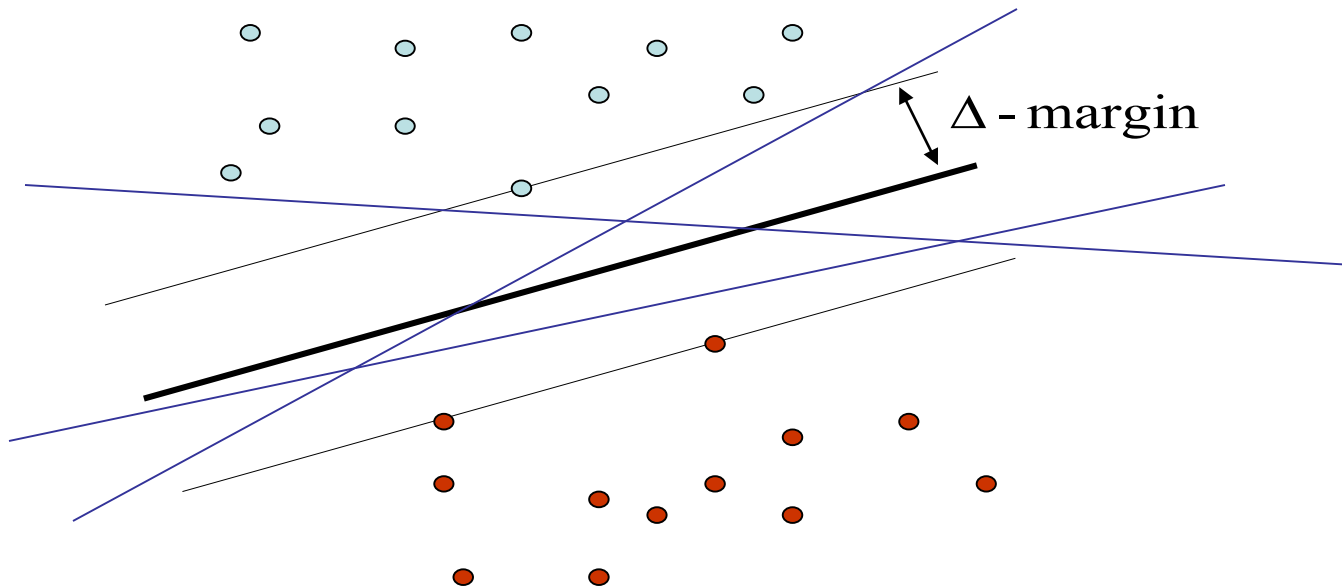
Support vector machine

Suppose that the data

$$(y_1, x_1), \dots, (y_l, x_l) \quad x \in \mathbb{R}^n \quad y \in \{+1, -1\}$$

can be separated by a hyperplane

$$(w \cdot x) - b = 0$$

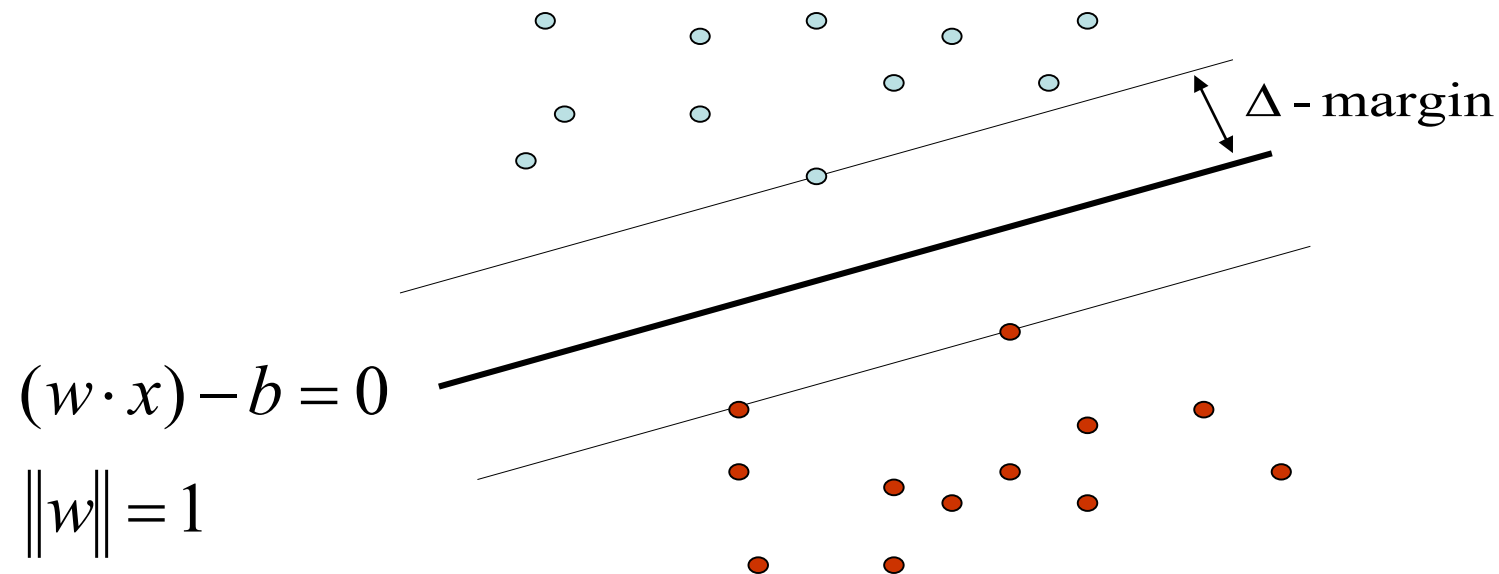


Support vector machine

Blue dots: $(w \cdot x_i) - b \geq 0, \quad y_i = +1$

Red dots: $(w \cdot x_i) - b \leq 0, \quad y_i = -1$

Combined: $y_i [(w \cdot x_i) - b] \geq 0, \quad \forall i$ Variables: w and b



Support vector machine

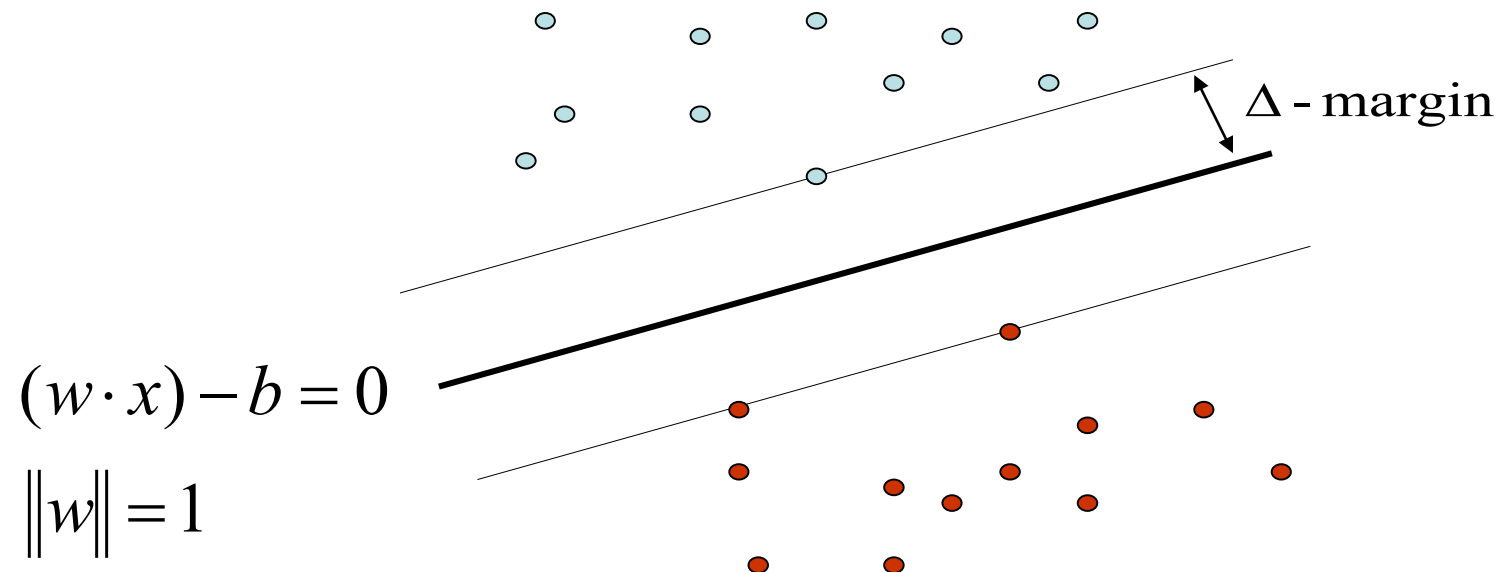
Blue dots: $(w \cdot x_i) - b \geq +\Delta, \quad y_i = +1$

$$\Delta \geq 0$$

Red dots: $(w \cdot x_i) - b \leq -\Delta, \quad y_i = -1$

Combined: $y_i[(w \cdot x_i) - b] \geq \Delta, \quad \forall i$

Variables: w and b

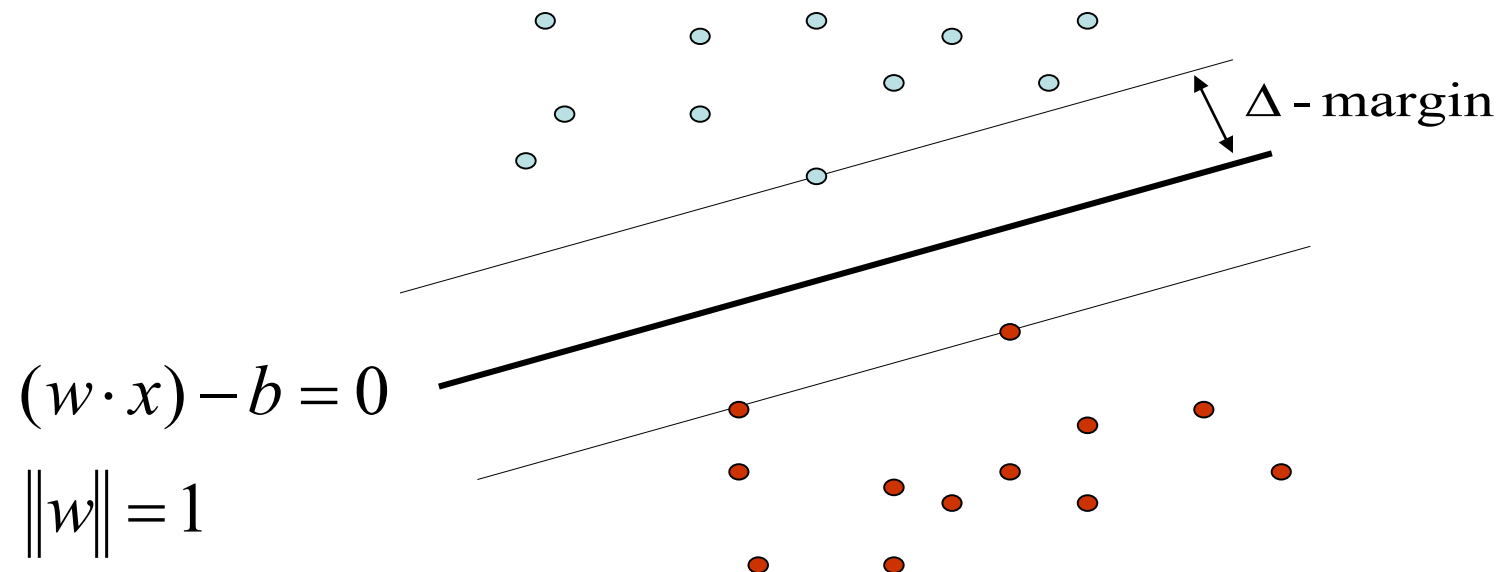


Support vector machine

Maximize the margin: $\max \Delta$

$$\text{s.t.} \quad y_i[(w \cdot x_i) - b] \geq \Delta, \quad \forall i$$
$$w_1^2 + \dots + w_n^2 = 1, \quad \Delta \geq 0$$

Variables: w , Δ and b



Support vector machine

$$y_i \left[\left(\frac{w}{\Delta} \cdot x_i \right) - \frac{b}{\Delta} \right] \geq 1, \quad \forall i \quad \text{or:} \quad y_i \left[(\bar{w} \cdot x_i) - \bar{b} \right] \geq 1, \quad \forall i$$

$$\bar{w} = w / \Delta, \quad \|\bar{w}\| = 1 / \Delta$$

Maximize the margin: $\min \|\bar{w}\|$

$$\text{s.t.} \quad y_i \left[(\bar{w} \cdot x_i) - \bar{b} \right] \geq 1, \quad \forall i$$

Variables: \bar{w} and \bar{b}

Support vector machine

Maximize the margin: $\min \|w\|^2$

$$\text{s.t. } y_i[(w \cdot x_i) - b] \geq 1, \quad \forall i$$

Variables: w and b

Support vector machine

Maximize the margin: $\min (w \cdot w)$

$$\text{s.t. } y_i[(w \cdot x_i) - b] \geq 1, \quad \forall i$$

Variables: w and b

Support vector machine

Maximize the margin: $\min 0.5(w \cdot w)$

$$\text{s.t. } y_i[(w \cdot x_i) - b] \geq 1, \quad \forall i$$

Variables: w and b

Non separable case:

Maximize the margin: $\min 0.5(w \cdot w) + C \left(\sum_{i=1}^l \xi_i \right)$

$$\text{s.t. } y_i[(w \cdot x_i) - b] \geq 1 - \xi_i, \quad \forall i$$

$$\xi_i \geq 0$$

Variables: w, b and ξ

Support vector machine

$$(y_1, x_1), \dots, (y_l, x_l) \quad x \in \mathbb{R}^n$$

Primal problem

$$\min 0.5(w \cdot w) + C \sum_{i=1}^l \xi_i$$

$$\text{s.t.} \quad \xi_i \geq 0$$

Variables: w, b and ξ

$$y_i [(x_i \cdot w) - b] \geq 1 - \xi_i, \quad i = 1, \dots, l$$

Dual problem

$$\max \sum_{i=1}^l \alpha_i - 0.5 \sum_{i,j}^l \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

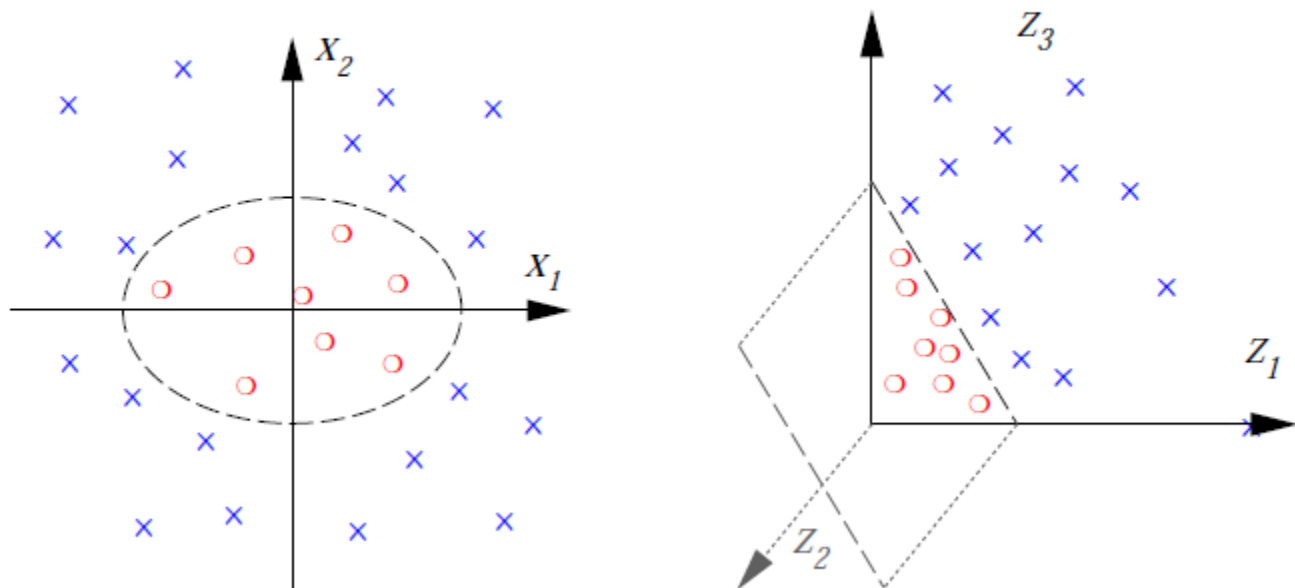
$$\text{s.t.}$$

Variables: α

$$\sum_{i=1}^l \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l$$

Kernels

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$



$$\begin{aligned}\langle \Phi(x), \Phi(x') \rangle &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2) (x_1'^2, \sqrt{2} x_1' x_2', x_2'^2)^\top \\ &= \langle x, x' \rangle^2 \\ &=: k(x, x')\end{aligned}$$

→ the dot product in \mathcal{H} can be computed in \mathbb{R}^2

Support vector machine

**Optimization
problem for finding
support vectors**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^l \alpha_i - 0.5 \sum_{i,j}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & \\ & \sum_{i=1}^l \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l \end{aligned}$$

Kernels

Polynomial machine:

$$K(x_i, x_j) = [\alpha(x_i \cdot x_j) + \beta]^d$$

A radial basis function machine:

$$K(x_i, x_j) = \exp \left\{ -\gamma \|x_i - x_j\|^2 \right\}$$

Support vector machine

**Optimization
problem for finding
support vectors**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^l \alpha_i - 0.5 \sum_{i,j}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & \\ & \sum_{i=1}^l \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l \end{aligned}$$

Decision rules with a kernel using found α_i^*

if $(\sum_{i=1}^l y_i \alpha_i^* K(x_i, x)) - b \geq 0$ then x is *blue*

if $(\sum_{i=1}^l y_i \alpha_i^* K(x_i, x)) - b < 0$ then x is *red*

$b = (\sum_{i=1}^l y_i \alpha_i^* K(x_i, x_{i_0})) - y_{i_0}$ for some $\alpha_{i_0}^* : 0 < \alpha_{i_0}^* < C, (\alpha_{i_0}^* \neq 0, \alpha_{i_0}^* \neq C)$

if there is no such $\alpha_{i_0}^*$, increase C and train again

Support vector machine

x_i that correspond to positive α_i^* are called the support vectors!!!

Only the support vectors carry important information!!!

They correspond to the active constraints of the primal problem!!!

Let $I^* = \{i : \alpha_i > 0\}$ be the set of support vectors

Decision rules using only the support vectors

if $(\sum_{i \in I^*} y_i \alpha_i^* K(x_i, x)) - b \geq 0$ then x is *blue*

if $(\sum_{i \in I^*} y_i \alpha_i^* K(x_i, x)) - b < 0$ then x is *red*

$b = (\sum_{i \in I^*} y_i \alpha_i^* K(x_i, x_{i_0})) - y_{i_0}$ for some $\alpha_{i_0}^* : 0 < \alpha_{i_0}^* < C, (\alpha_{i_0}^* \neq 0, \alpha_{i_0}^* \neq C)$

if there is no such $\alpha_{i_0}^*$, increase C and train again

Support vector machine

**Optimization
problem for finding
support vectors**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^l \alpha_i - 0.5 \sum_{i,j}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ \text{s.t.} \quad & \\ & \sum_{i=1}^l \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l \end{aligned}$$

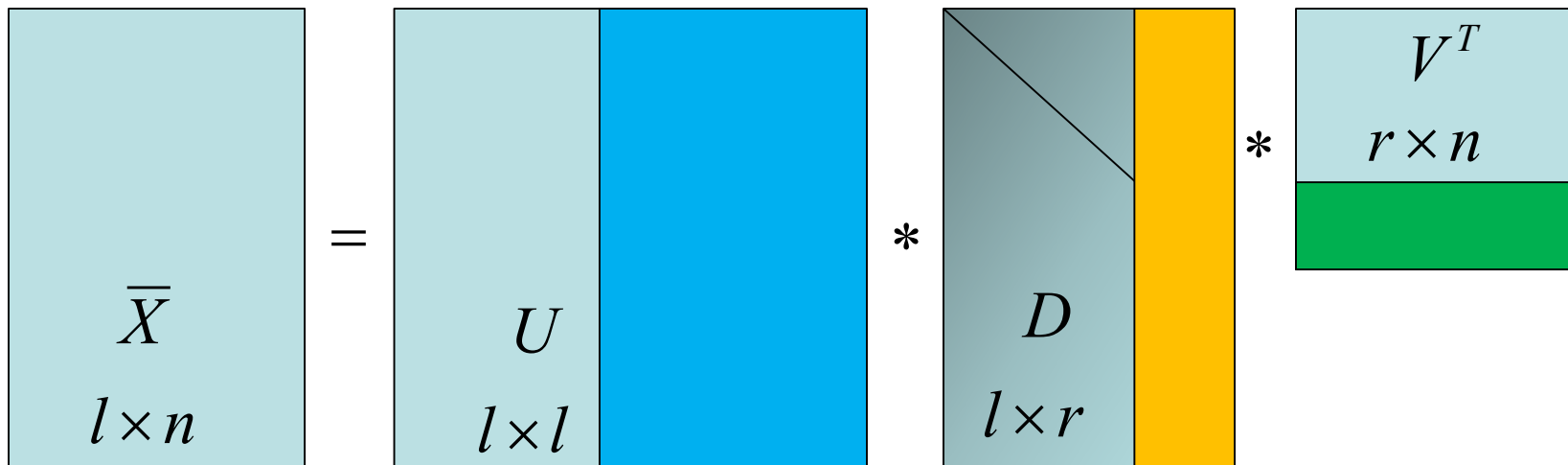
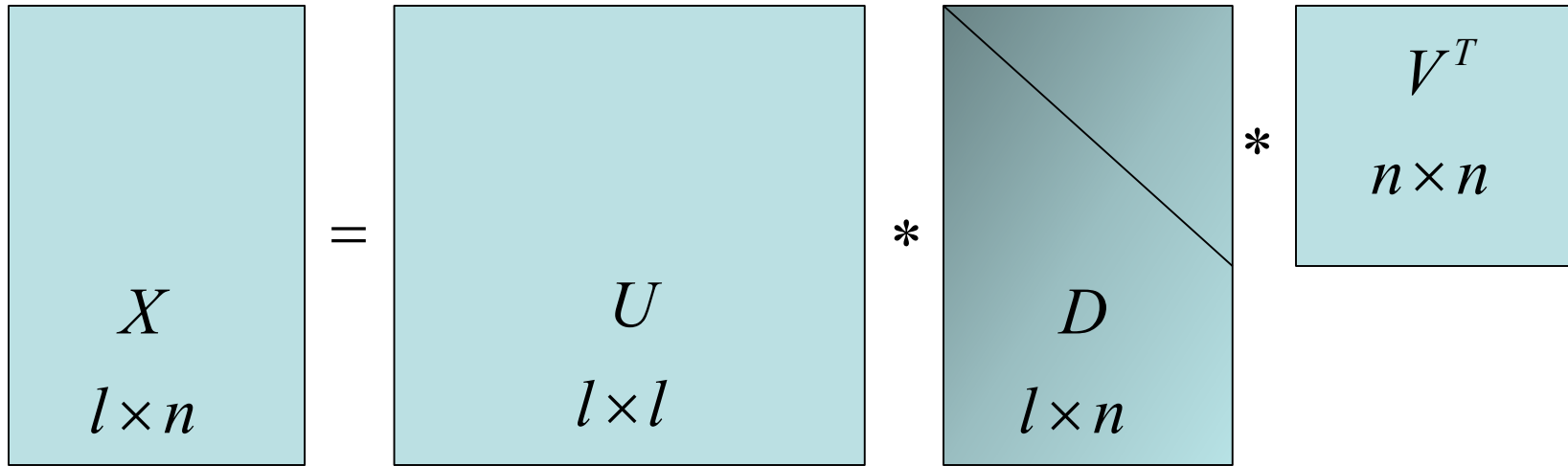
Matlab QP setting

$$\begin{aligned} \min_{\alpha} \quad & 0.5 \alpha^T M \alpha - e^T \alpha \\ \text{s.t.} \quad & \\ & y^T \alpha = 0, \quad 0 \leq \alpha \leq C e \\ & \text{where } M_{ij} = y_i y_j K(x_i, x_j), \\ & y = (y_1, \dots, y_l)^T, e = (1, \dots, 1)^T \end{aligned}$$

Linear Principle Component Analysis = Singular Value Decomposition of X

$$X = UDV^T$$

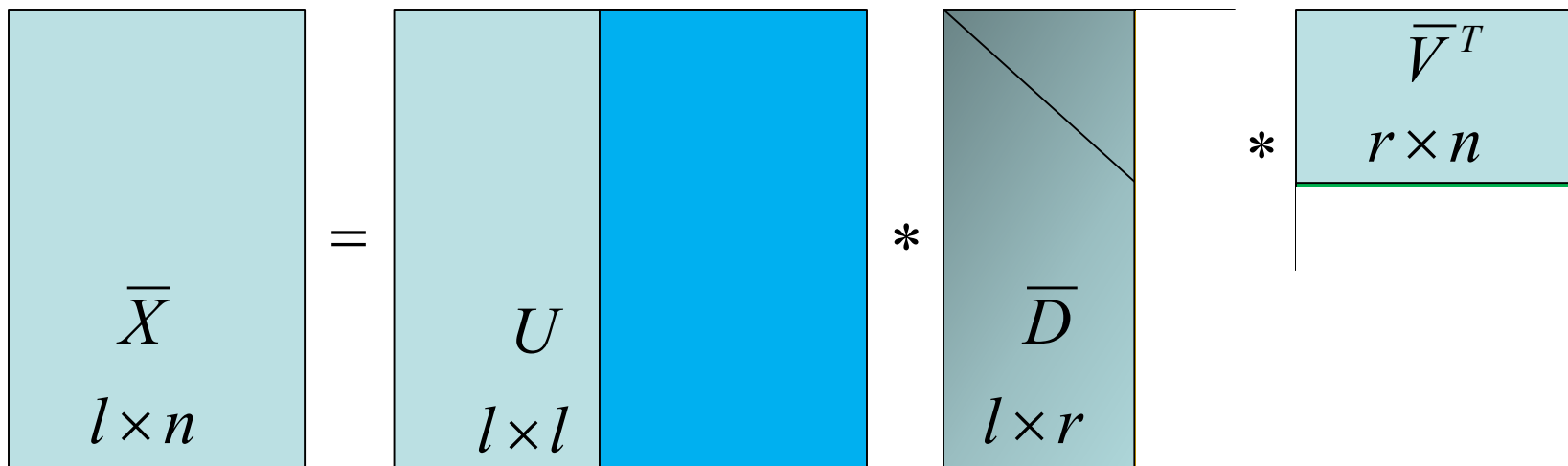
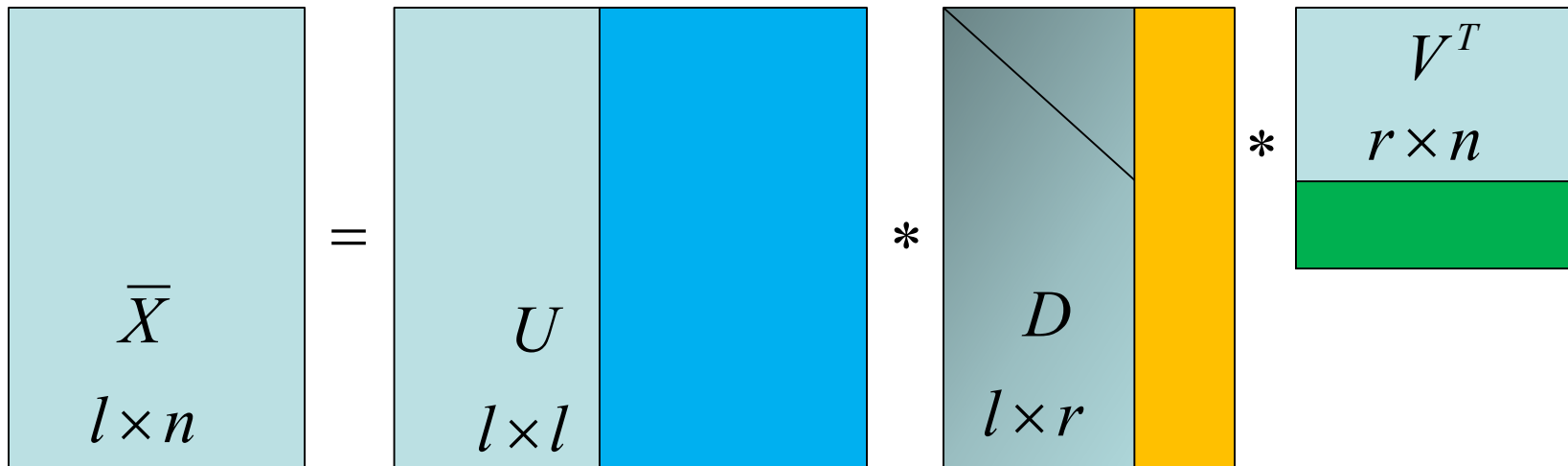
$$XV = UD$$



Linear Principle Component Analysis = Singular Value Decomposition of X

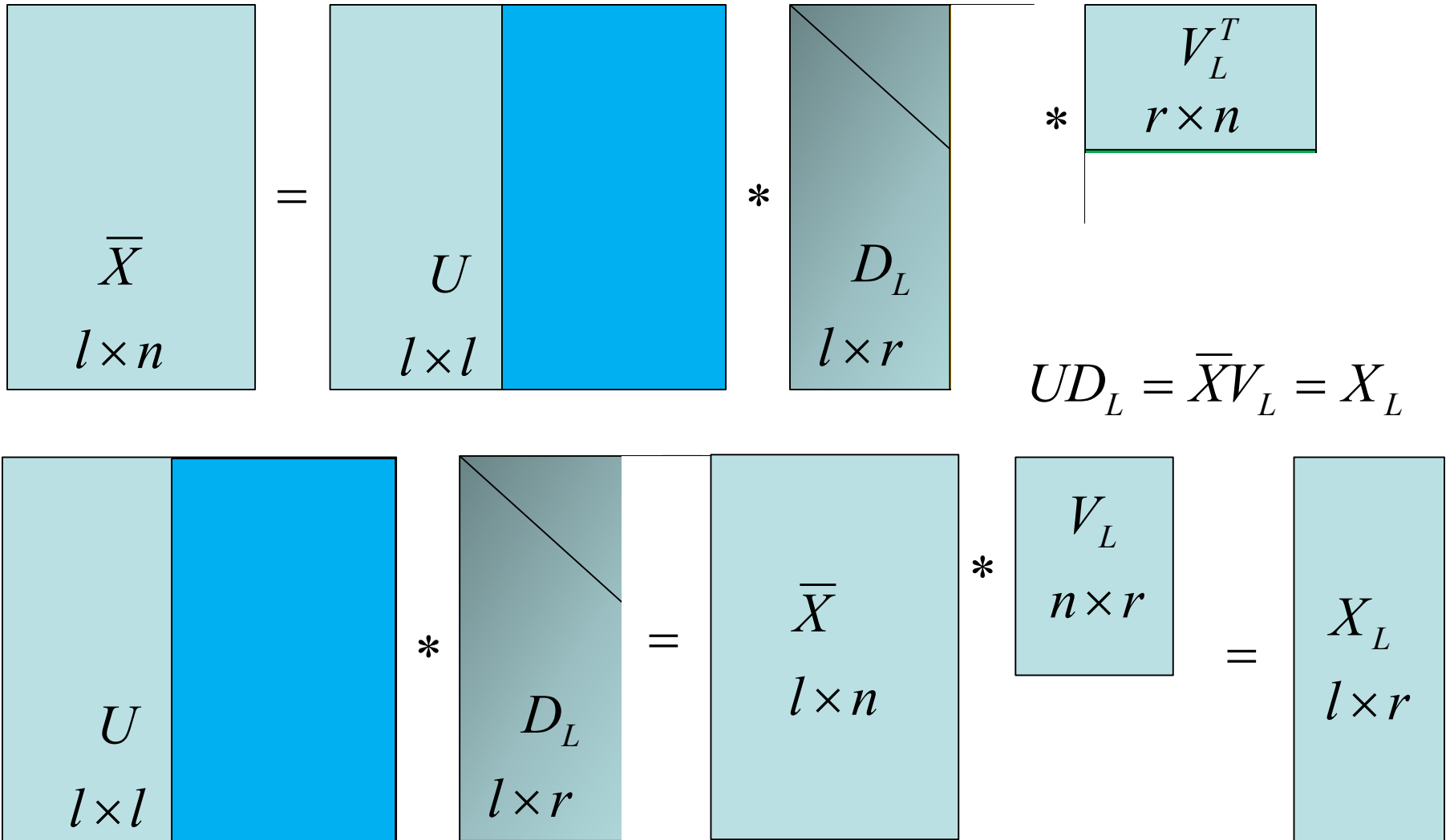
$$\bar{X} = U D_L V_L^T$$

$$\bar{X} V_L = U D_L = X_L$$



Linear Principle Component Analysis = Singular Value Decomposition of X

$$\bar{X} = UD_L V_L^T$$



SVM testing with PCA

1. Calculate the SVD: $X = UDV^T$

2. Reduce the dimensionality of the feature space:

$$X_L^{training} = UD_L, \text{ for training data}$$

$$X_L^{testing} = X^{testing} V_L, \text{ for testing data,}$$

V_L is calculated on the training data

3. Perform the SVM as usual