

Machine Learning Algorithms

Optimization Problems for Machine Learning

Markus Rupp

13.9.2020



Problem Formulation

- Given two sets S_1 and S_2 with presence and absence of feature A, respectively, we try to find a minimum (maximum)

$$\min_{\underline{\hat{w}}} \sum_{\underline{x}_i \in S_1} \left(1 - \underline{\hat{w}}^H \underline{\hat{x}}_i\right)^2 + \sum_{\underline{x}_i \in S_2} \left(-1 - \underline{\hat{w}}^H \underline{\hat{x}}_i\right)^2$$

- As w is underdetermined, we also like to constrain it
- For example the smallest, which can be achieved by adding an additional constraint

$$\min_{\underline{\hat{w}}} \sum_{\underline{x}_i \in S_1} \left(1 - \underline{\hat{w}}^H \underline{\hat{x}}_i\right)^2 + \sum_{\underline{x}_i \in S_2} \left(-1 - \underline{\hat{w}}^H \underline{\hat{x}}_i\right)^2 + \lambda \left\| \underline{\hat{w}} \right\|^2$$

Problem Formulation

- This has an additional advantage, as often saddle points occur when nonlinear formulations $\sigma(x)$ are included, e.g.,

$$\min_{\underline{\tilde{w}}} \sum_{i \in S_1} \left(1 - \sigma \left(\underline{\tilde{w}}^H \underline{\tilde{x}}_i \right) \right)^2 + \sum_{i \in S_2} \left(-1 - \sigma \left(\underline{\tilde{w}}^H \underline{\tilde{x}}_i \right) \right)^2 + \lambda \left\| \underline{\tilde{w}} \right\|^2$$

- With the additional term the saddle point is often deformed into a continuous monotone region which can be searched through with a gradient approach.

Interpretation

- The additional term can be interpreted in different ways
- 1) additive weight with $\lambda > 0$
- 2) Lagrangian Multiplier
- 3) reformulated problem with side constraints



Interpretation 1: additive weight

- In a linear classification setting, adding a positive term changes the LS solution :

$$R = \sum_{\hat{\underline{x}}_i \in S} \hat{\underline{x}}_i \hat{\underline{x}}_i^H \rightarrow R = \lambda I + \sum_{\hat{\underline{x}}_i \in S} \hat{\underline{x}}_i \hat{\underline{x}}_i^H$$

- Regularisation term to ensure positive definiteness of R.
A very small $\lambda > 0$ will do!

Interpretation 1: additive weight

- Consequence for the gradient approach:

$$\underline{\hat{w}}_k = (1 - \lambda) \underline{\hat{w}}_{k-1} + \mu \underline{x}_i (y_i - \underline{\hat{w}}_{k-1}^H \underline{x}_i)$$

- LMS obtains a „leaky“ term,
forgetting a bit of its previous estimate



Interpretation 2: Lagrangian Multiplier

- In this case λ is interpreted as Lagrangian multiplier and thus the minimum of \underline{w} becomes a constraint.
- Under all minimal \underline{w} we search for those that minimize the LS approach.

$$\frac{\partial}{\partial \underline{w}} \sum_{\hat{x}_i \in S_1} \left(1 - \underline{w}^H \hat{x}_i\right)^2 + \sum_{\hat{x}_i \in S_2} \left(-1 - \underline{w}^H \hat{x}_i\right)^2 + \lambda \left\| \underline{w} \right\|^2 = 0$$

$$-2\underline{p} + 2R\underline{w} + 2\lambda\underline{w} = 0$$

$$(R + \lambda I)\underline{w} = \underline{p}$$

$$\min_{\lambda} \sum_{\hat{x}_i \in S_1} \left(1 - \underline{p}^H (R + \lambda I)^{-1} \hat{x}_i\right)^2 + \sum_{\hat{x}_i \in S_2} \left(-1 - \underline{p}^H (R + \lambda I)^{-1} \hat{x}_i\right)^2 + \lambda \underline{p}^H (R + \lambda I)^{-2} \underline{p}$$

Nonlinear in λ
Can be solved
numerically



Interpretation 3: Reformulation

$$\min \|\underline{\widehat{w}}\|^2$$

subject to

$$\underline{\widehat{w}}^H \underline{\widehat{x}}_i = y_i; \quad i = 1, 2, \dots$$

$$\min \frac{1}{2} \|\underline{\widehat{w}}\|^2 + \sum \lambda_i \left(\underline{\widehat{w}}^H \underline{\widehat{x}}_i - y_i \right)$$

- Which turns out to be an other form of a Lagrangian optimization problem

Interpretation 3: Reformulation

$$\frac{\partial}{\partial \underline{\widehat{w}}} \left\{ \frac{1}{2} \|\underline{\widehat{w}}\|^2 + \sum \lambda_i \left(\underline{\widehat{w}}^H \underline{\widehat{x}}_i - y_i \right) \right\} = \underline{0}$$

$$\underline{\widehat{w}} + \sum \lambda_i \underline{\widehat{x}}_i = \underline{0}$$

$$\min \frac{1}{2} \left\| \sum \lambda_i \underline{\widehat{x}}_i \right\|^2 - \left\| \sum \lambda_i \underline{\widehat{x}}_i \right\|^2 - \sum \lambda_i y_i$$

$$= \min \frac{1}{2} \left\| \sum \lambda_i \underline{\widehat{x}}_i \right\|^2 + \sum \lambda_i y_i$$

Interpretation 3: Reformulation

- Which can be explicitly solved for λ_i : $X = [\underline{\hat{x}}_1, \underline{\hat{x}}_2, \dots, \underline{\hat{x}}_N] \rightarrow R = X^H X$
- Consider Gramian R built from $\underline{\hat{x}}_i$

$$R \underline{\lambda} = -\underline{y} \rightarrow \underline{\lambda} = -R^{-1} \underline{y}$$

$$\frac{1}{2} \|\underline{\hat{w}}\|^2 + \sum \lambda_i \left(\underline{\hat{w}}^H \underline{\hat{x}}_i - y_i \right)$$

$$= \frac{1}{2} \|\underline{\hat{w}}\|^2 + \underline{\lambda}^H \left(X \underline{\hat{w}} - \underline{y} \right)$$

$$= \frac{1}{2} \|\underline{\hat{w}}\|^2 - \underline{y}^H \left[X^H X \right]^{-1} \left(X \underline{\hat{w}} - \underline{y} \right)$$



Interpretation 3: Reformulation

- Now solving with respect to $\underline{\hat{w}}$ is simply

$$\min_{\underline{\hat{w}}} \frac{1}{2} \|\underline{\hat{w}}\|^2 - \underline{y} [X^H X]^{-1} (X \underline{\hat{w}} - \underline{y})$$

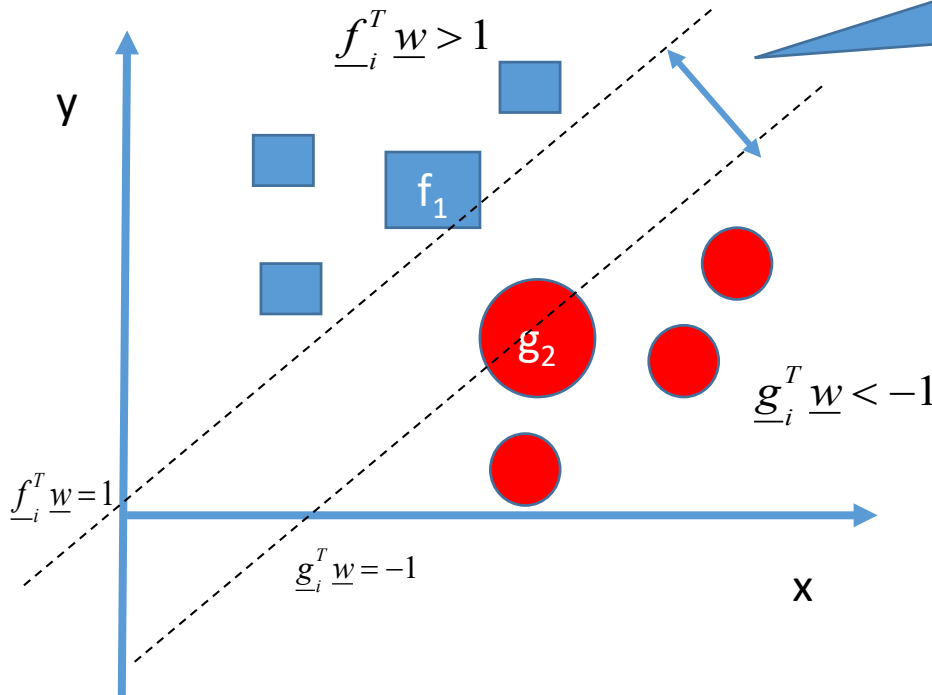
$$\underline{\hat{w}} = X^H [X^H X]^{-1} \underline{y}$$

Note that this is an underdetermined LS solution which would only work perfectly if the observations are less than parameters, which is **not** the case!



We need a better Re-Formulation: Let's go deeper (Vapnik)

- We do not want any hyperplane that separates two classes but an optimal line that guarantees that the border members of the classes are at maximum distance.



Let us maximize this distance here

$$\langle \underline{w}, \underline{f}_1 - \underline{g}_2 \rangle = 1 - (-1) = 2$$

$$\left\langle \frac{\underline{w}}{\|\underline{w}\|}, \underline{f}_1 - \underline{g}_2 \right\rangle = \frac{2}{\|\underline{w}\|}$$

which means
minimizing $\|\underline{w}\|$

\underline{f}_1 and \underline{g}_2 are supporting these separation lines
→ Vector support machine



We need a better Re-Formulation: Let's go deeper (Vapnik)

- Let us do this again but with more subtleties

- Recall:
$$S_1 : \quad \underline{\widehat{w}}^H \underline{\widehat{x}}_i \geq 1 \quad \rightarrow y_i = 1$$
$$\underline{w}^H \underline{x}_i + b \geq 1$$
$$y_i \left(\underline{w}^H \underline{x}_i + b \right) \geq 1$$
$$S_2 : \quad \underline{\widehat{w}}^H \underline{\widehat{x}}_i \leq -1 \quad \rightarrow y_i = -1$$
$$\underline{w}^H \underline{x}_i + b \leq -1$$
$$y_i \left(\underline{w}^H \underline{x}_i + b \right) \geq 1$$

Vector Support Machine

- Let's differentiate

Note the
difference

and here

$$\min_{\underline{\hat{w}}} \left\{ \frac{1}{2} \|\underline{\hat{w}}\|^2 - \sum \lambda_i \left(y_i \left(\underline{\hat{w}}^H \underline{\hat{x}}_i \right) - 1 \right) \right\}$$

$$= \min_{\underline{w}, b} \left\{ \frac{1}{2} \|\underline{w}\|^2 - \sum \lambda_i \left(\left(\underline{w}^H \underline{x}_i + b \right) y_i - 1 \right) \right\}$$

$$\frac{\partial}{\partial \underline{w}, b} \left\{ \frac{1}{2} \|\underline{w}\|^2 - \sum \lambda_i \left(\left(\underline{w}^H \underline{x}_i + b \right) y_i - 1 \right) \right\} = \underline{0}$$



Vector Support Machine (Vapnik 1963)

- Let's differentiate

$$\frac{\partial}{\partial \underline{w}} \left\{ \frac{1}{2} \|\underline{w}\|^2 - \sum \lambda_i \left((\underline{w}^H \underline{x}_i + b) y_i - 1 \right) \right\} = \underline{0}$$
$$\underline{w} = \sum \lambda_i y_i \underline{x}_i$$

$$\frac{\partial}{\partial b} \left\{ \frac{1}{2} \|\underline{w}\|^2 - \sum \lambda_i \left((\underline{w}^H \underline{x}_i + b) y_i - 1 \right) \right\} = 0$$
$$0 = \sum \lambda_i y_i$$

solution by minimizing this
w.r.t. λ_i :

$$-\frac{1}{2} \left\| \sum \lambda_i y_i \underline{x}_i \right\|^2 = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \underline{x}_i, \underline{x}_j \rangle$$

Quadratic in λ ,
thus can be
minimized.

Note that the
minimum of the
cost function
depends only on
the dot product

Vector Support Machine

- Given the training data $\{\underline{x}_i, y_i\}$

$$\min_{\lambda_j} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \underline{x}_i, \underline{x}_j \rangle$$

- By maximizing the dual problem:

$$\max_{\lambda_j} \sum_j \lambda_j - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \underline{x}_i, \underline{x}_j \rangle$$

Only support vectors
have $\lambda_i > 0$

- That satisfies

$$\lambda_i \geq 0 \quad \wedge \quad \sum \lambda_i y_i = 0$$

Vector Support Machine

- With the support vectors compute

$$\underline{w} = \sum \lambda_i y_i \underline{x}_i$$

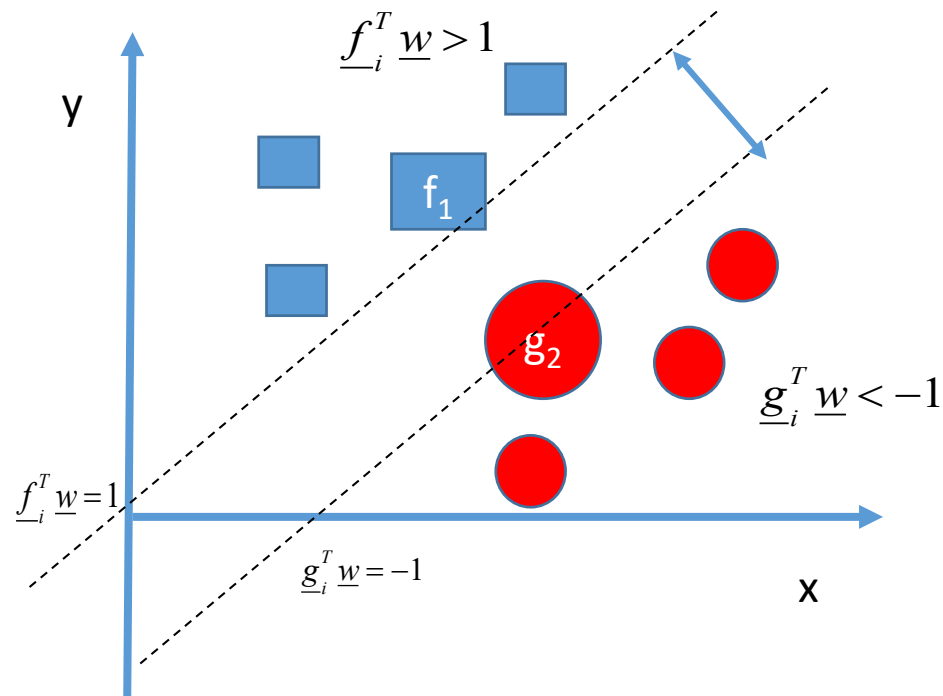
- Select support vector set $\{\underline{f}_1, y_1=1\}$

$$\underline{w}^T \underline{f}_1 + b = 1$$

$$b = 1 - \underline{w}^T \underline{f}_1$$

Vector Support Machine Example

- Recall our initial problem



We have two support vectors \underline{f}_1 and \underline{g}_2 with $y=1$ and -1 , respectively
This results in $\lambda_1=\lambda_2$

We need to

$$\max_{\lambda_j} \sum_j \lambda_j - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \underline{x}_i, \underline{x}_j \rangle$$

$$= \max 2\lambda_1 - \frac{1}{2} \lambda_1^2 \|\underline{f}_1 - \underline{g}_2\|_2^2$$

$$\lambda_1 = \lambda_2 = \frac{2}{\|\underline{f}_1 - \underline{g}_2\|_2^2}$$

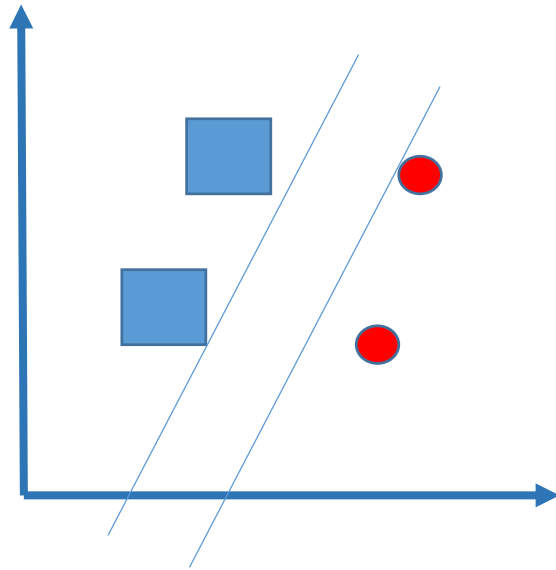
$$\underline{w} = \frac{2}{\|\underline{f}_1 - \underline{g}_2\|_2^2} (\underline{f}_1 - \underline{g}_2)$$

$$b = 1 - \frac{2}{\|\underline{f}_1 - \underline{g}_2\|_2^2} (\underline{f}_1 - \underline{g}_2)^T \underline{f}_1$$

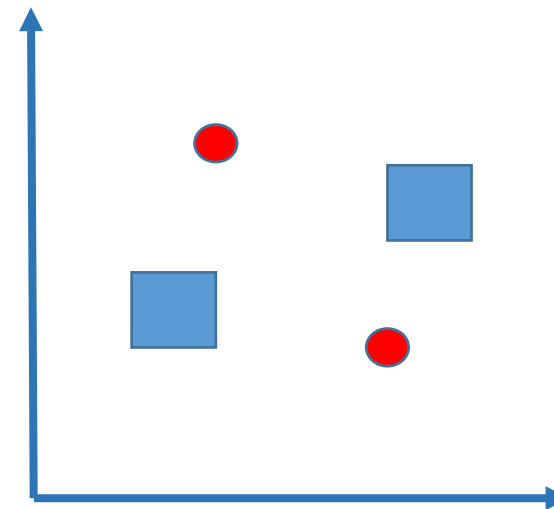


Recall Linear Separable

- Linearly separable



non-linearly separable



Classical: the XOR problem



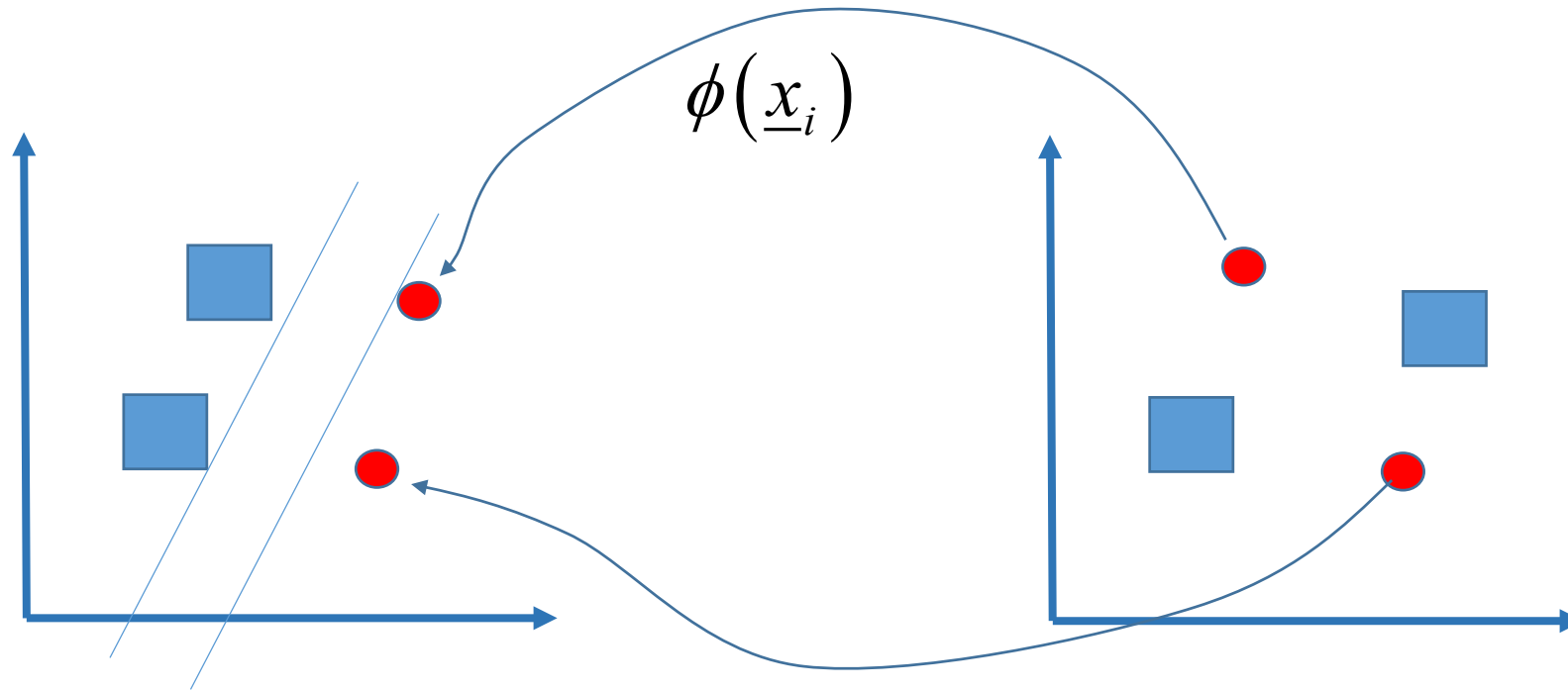
Ways to obtain linearly separable sets

- Nonlinear mapping of input vectors
- Increasing dimensions of the feature space

Recall Linear Separable nonlinear mapping of input vectors

- Linearly separable

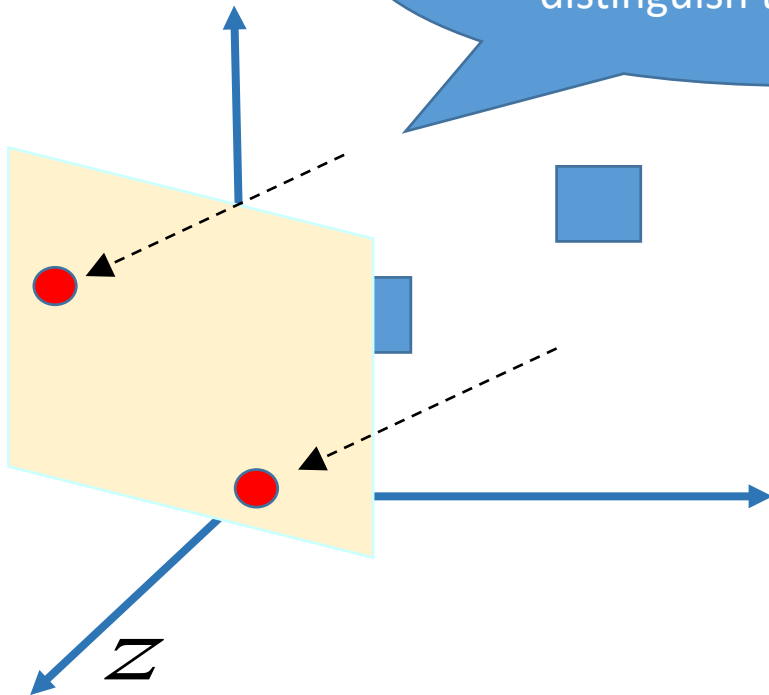
non-linearly separable



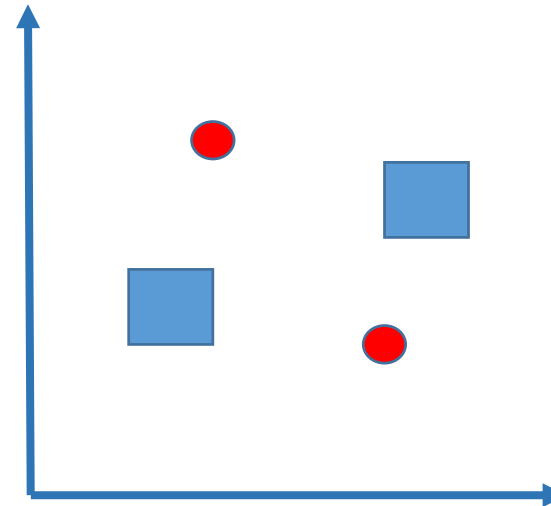
Recall Linear Separable

- Linearly separable

Moving red circles to higher z allows a hyperplane to distinguish the two classes



non-linearly separable



Moving into higher dimensions (Vapnik 1992)

- For example use a nonlinear mapping of input values

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \phi[\underline{x}] = \begin{bmatrix} \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

$$\left(1 + \langle \underline{x}, \underline{y} \rangle\right)^2 - 1 = \langle \phi[\underline{x}], \phi[\underline{y}] \rangle$$

- and obtain a polynomial kernel

Now back to our original problem

$$\begin{aligned} L &= -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \underline{x}_i, \underline{x}_j \rangle \\ &= -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j K(\underline{x}_i, \underline{x}_j) \end{aligned}$$

- Interpret dot product as a kernel
- Kernel defines the (non)-linear mapping

$$K(\underline{x}, \underline{y}) = \langle \phi(\underline{x}), \phi(\underline{y}) \rangle$$



Kernels

Vladimir Naumovich Vapnik ([Russian](#): Владимир Наумович Вапник;
born 6 December 1936)
~1990 at AT&T Bell Labs

- $K_0(\underline{x}, \underline{y}) = \tanh(\beta^2 \underline{x}^H \underline{y} + \alpha)$ Neuronal Network
- $K_1(\underline{x}, \underline{y}) = (\underline{x}^H \underline{y} + 1)^n - 1$ Polynomial Kernel of degree n
- $K_2(\underline{x}, \underline{y}) = e^{-\beta \|\underline{x} - \underline{y}\|^2}$ Radial Basis Function (RBF)
- $K_3(\underline{x}, \underline{y}) = \prod_{n=1}^N \frac{\sin((2D+1)\pi(x_n - y_n))}{\sin(\pi(x_n - y_n))} - 1$ Fourier

