



AARHUS
UNIVERSITET

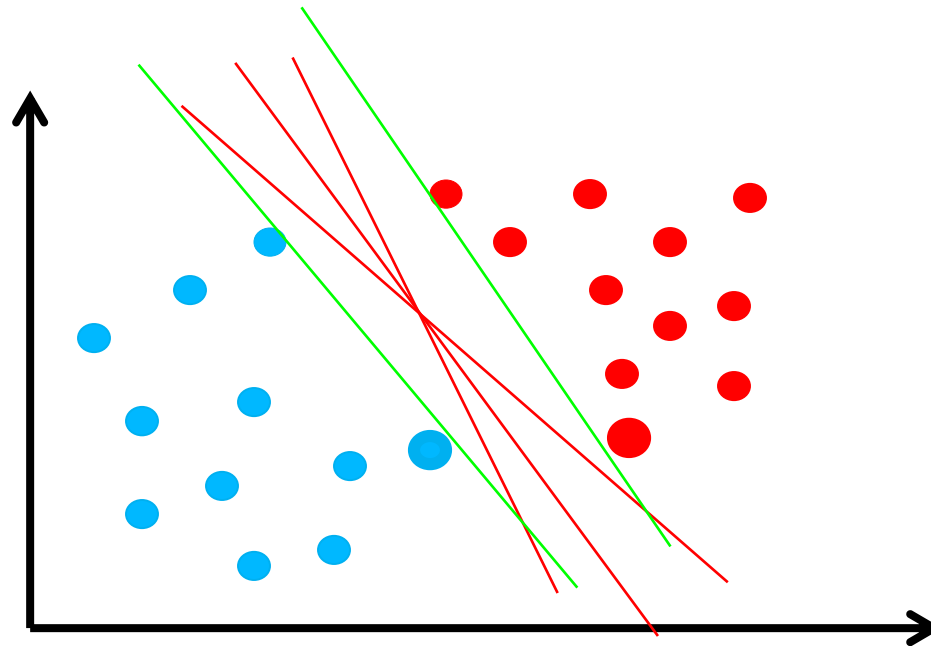
Optimization and Data Analytics

Alexandros Iosifidis

@

Aarhus University, Department of Engineering

Support Vector Machine



Infinite decision functions possible

Support Vector Machine

That is, given a set of N samples, each represented by a vector $\mathbf{x}_i \in \mathbb{R}^D$, and the corresponding labels $l_i = \{-1, 1\}$ we want to optimize the parameters of $g(\cdot)$ in order to define a discriminant hyperplane discriminating the two classes.

Support Vector Machine (SVM) assumes the data \mathbf{x} is mapped to $\boldsymbol{\phi}$ using a function $\phi(\cdot)$

$$\mathbf{x}_i \in \mathbb{R}^D \xrightarrow{\phi(\cdot)} \boldsymbol{\phi}_i \in \mathcal{F}$$

Note that the above is a generic mapping. For example a linear function $\phi(\mathbf{x}) = \mathbf{x}$ can also be used.

Then, we define the decision function

$$g(\boldsymbol{\phi}_i) = \mathbf{w}^T \boldsymbol{\phi}_i - b$$

Support Vector Machine

If the parameters of the decision function are optimized, then

$$l_i g(\phi_i) \geq 0 \Rightarrow l_i (\mathbf{w}^T \phi_i - b) \geq 0$$

or

$$\begin{aligned} \mathbf{w}^T \phi_i - b &\geq q, \text{ for } l_i = 1 \text{ and} \\ \mathbf{w}^T \phi_i - b &\leq -q, \text{ for } l_i = -1. \end{aligned}$$

q expresses the minimal distance between the decision hyperplane and the closest to it training samples. That is, q is the margin appearing between the two classes.

Notice that there are multiple w 's (see next slide).

But we only care about when w is the unit vector

Remember that the decision function expresses distance of a sample from the hyper-plane.

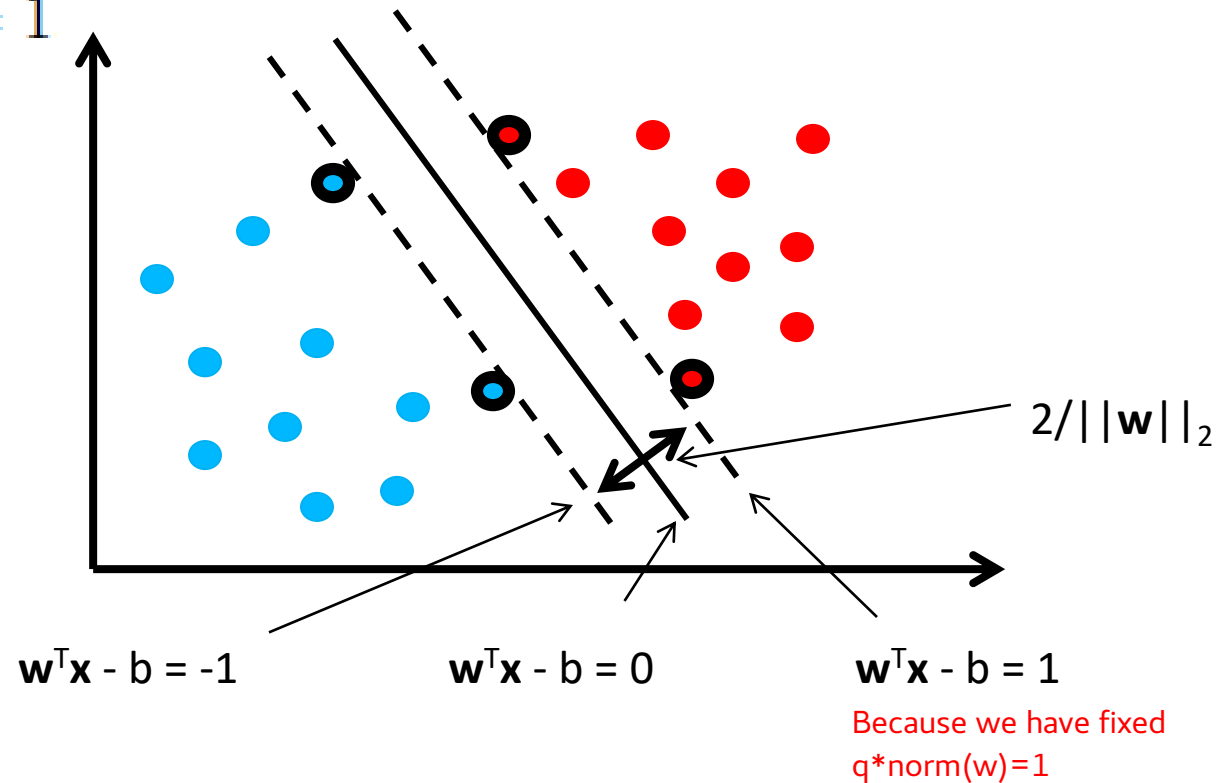
Thus for ϕ_m
(the closest point)

$$\frac{l_m g(\phi_m)}{\|\mathbf{w}\|_2} = q$$

Support Vector Machine

We fix the right side of the equation in the previous slide

We use: $q||\mathbf{w}||_2 = 1$



Support Vector Machine

In order to define the weights w and the margin b , SVM optimizes for

Minimise the inverse of w which maximises the margin

$$\mathcal{J}_{SVM} = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i,$$

subject to the constraints:

$$\begin{aligned} l_i(\mathbf{w}^T \phi_i - b) &\geq 1 - \xi_i, \quad i = 1, \dots, N \\ \xi_i &\geq 0. \end{aligned}$$

When the C_i is zero
then the corresponding x
is in the margin.

Between 0 to 1, it is on the
correct side of the decision
function but inside the margin.

If $C_i > 1$ then the sample will
be on the wrong side of the
decision function

Support Vector Machine

To optimize J_{SVM} s.t. the constraints, we define the Lagrangian

$$\mathcal{L} = \frac{1}{2} \mathbf{w}^T \mathbf{w} + c \sum_{i=1}^N \xi_i - \sum_{i=1}^N \beta_i \xi_i - \sum_{i=1}^N \alpha_i [l_i(\mathbf{w}^T \phi_i - b) - 1 + \xi_i]$$

It is possible to use a fixed alpha and beta instead of alphas and betas for each sample.

Support Vector Machine

To optimize J_{SVM} s.t. the constraints, we define the Lagrangian

$$\mathcal{L} = \frac{1}{2} \mathbf{w}^T \mathbf{w} + c \sum_{i=1}^N \xi_i - \sum_{i=1}^N \beta_i \xi_i - \sum_{i=1}^N \alpha_i [l_i (\mathbf{w}^T \phi_i - b) - 1 + \xi_i]$$

The derivatives of \mathcal{L} w.r.t. all variables are

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 &\Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i l_i \phi_i, \\ \frac{\partial \mathcal{L}}{\partial b} = 0 &\Rightarrow \sum_{i=1}^N \alpha_i l_i = 0, \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = 0 &\Rightarrow c - \alpha_i - \beta_i = 0. \end{aligned}$$

Support Vector Machine

If we substitute these equations to L, we obtain

We compute the results and get a quadratic function

$$\max_{\alpha} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j l_i l_j \phi_i^T \phi_j + \sum_{i=1}^N \alpha_i$$

subject to the constraints

$$0 \leq \alpha_i \leq c, i = 1, \dots, N$$

Support Vector Machine

If we substitute these equations to L, we obtain

$$\max_{\alpha} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j l_i l_j \phi_i^T \phi_j + \sum_{i=1}^N \alpha_i$$

subject to the constraints

$$0 \leq \alpha_i \leq c, i = 1, \dots, N$$

Which can also be written as

$$\max_{\alpha} \alpha^T (\mathbf{l}^T \circ \mathbf{K}) \alpha + \mathbf{1}^T \alpha$$

where $\mathbf{K} = \Phi^T \Phi$

Support Vector Machine

The problem

$$\max_{\alpha} \alpha^T (\mathbf{l}^T \circ \mathbf{K}) \alpha + \mathbf{1}^T \alpha$$

is a quadratic problem having one global solution when \mathbf{K} is positive semi-definite.

After obtaining α , \mathbf{w} is calculated by
$$\mathbf{w} = \sum_{i=1}^N \alpha_i l_i \phi_i$$

b can be calculated by selecting a training sample for which $\alpha_i > 0$ and computing

$$b = \mathbf{w}^T \phi_i - l_i$$

Kernels

We can use any function $\kappa(+, +)$ defined on vector-pairs in order to calculate the elements of a matrix \mathbf{K}_{ij} as long as the resulting matrix \mathbf{K} is positive semi-definite.

Some example functions are

$$\begin{aligned}\kappa(\mathbf{x}_i, \mathbf{x}_j) &= e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2} && \text{There is an equivalent } \sigma = 1/(2\sigma^2) \\ \kappa(\mathbf{x}_i, \mathbf{x}_j) &= (\mathbf{x}_i^T \mathbf{x}_j + 1)^d\end{aligned}$$

Using a generic function $\kappa(+, +)$, we have

$$g(\mathbf{x}_*) = \mathbf{w}^T \phi_* - b = \sum_{i=1}^N l_i \alpha_i \phi_i^T \phi_* - b = \boldsymbol{\alpha}^T \mathbf{L} \mathbf{k}_* - b$$

Kernel Least-Means Square Regression

We can also define a linear regression using the data $\phi_i \in \mathcal{F}$, $i = 1, \dots, N$

$$\mathbf{W}^T \phi_i = \mathbf{t}_i, \quad i = 1, \dots, N$$

In order to determine the matrix \mathbf{W} , we optimize for

$$\begin{aligned} J_{LSE} &= \|\mathbf{W}^T \Phi - \mathbf{T}\|_F^2 \\ &= \text{Tr}(\mathbf{W}^T \Phi \Phi^T \mathbf{W} - 2\mathbf{W}^T \Phi \mathbf{T} + \mathbf{T} \mathbf{T}^T) \end{aligned}$$

where $\Phi = [\phi_1, \dots, \phi_N]$, $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_N]$ and $\text{Tr}(\cdot)$ is the trace operator of a matrix.

Kernel Least-Means Square Regression

Setting the derivative w.r.t. \mathbf{W} equal to zero, we have

$$\nabla \mathcal{J}_{LSE} = 0 \Rightarrow 2\Phi\Phi^T\mathbf{W} = 2\Phi\mathbf{T}^T$$

leading to

$$\mathbf{W} = (\Phi\Phi^T)^{-1} \Phi\mathbf{T}^T = \Phi^\dagger\mathbf{T}^T$$

when the mapping $\mathbf{x} \rightarrow \boldsymbol{\phi}$ is defined, we can use the above equation to calculate \mathbf{W} .

Kernel Least-Means Square Regression

When the mapping $\mathbf{x} \rightarrow \boldsymbol{\phi}$ is defined through the function $\kappa(+,+)$, we express \mathbf{W} as a linear combination of the training samples

$$\mathbf{W} = \Phi \mathbf{A}$$

Representer theorem

Substituting \mathbf{W} to J_{LSE} , we obtain

$$\begin{aligned} \mathcal{J}_{LSE} &= \|\mathbf{A}^T \Phi^T \Phi - \mathbf{T}\|_F^2 = \|\mathbf{A}^T \mathbf{K} - \mathbf{T}\|_F^2 \\ &= \text{Tr}(\mathbf{A}^T \mathbf{K} \mathbf{K}^T \mathbf{A} - 2\mathbf{A}^T \mathbf{K} \mathbf{T} + \mathbf{T} \mathbf{T}^T) \end{aligned}$$

Why does \mathbf{K} need to be positive semi-definite? Suppose we have $\mathbf{K} = \mathbf{U} \mathbf{S} \mathbf{U}^T$, then we can write $\mathbf{U} \mathbf{S}^{(1/2)} \mathbf{U}^T \mathbf{U}^T \mathbf{S}^{(1/2)} \mathbf{U}^T$

Kernel Least-Means Square Regression

When the mapping $\mathbf{x} \rightarrow \boldsymbol{\phi}$ is defined through the function $\kappa(+,+)$, we express \mathbf{W} as a linear combination of the training samples

$$\mathbf{W} = \Phi \mathbf{A}$$

Substituting \mathbf{W} to \mathcal{J}_{LSE} , we obtain

$$\begin{aligned}\mathcal{J}_{LSE} &= \|\mathbf{A}^T \Phi^T \Phi - \mathbf{T}\|_F^2 = \|\mathbf{A}^T \mathbf{K} - \mathbf{T}\|_F^2 \\ &= \text{Tr}(\mathbf{A}^T \mathbf{K} \mathbf{K}^T \mathbf{A} - 2\mathbf{A}^T \mathbf{K} \mathbf{T} + \mathbf{T} \mathbf{T}^T)\end{aligned}$$

Setting the derivative to zero we have $\nabla \mathcal{J}_{LSE} = 0 \Rightarrow 2\mathbf{K} \mathbf{K}^T \mathbf{A} = 2\mathbf{K} \mathbf{T}^T$

$$\mathbf{A} = (\mathbf{K} \mathbf{K}^T)^{-1} \mathbf{K} \mathbf{T}^T = \mathbf{K}^\dagger \mathbf{T}^T$$

Kernel Discriminant Analysis

We can also define a linear projection using the data $\phi_i \in \mathcal{F}$, $i = 1, \dots, N$

In this case, the two scatter matrices are

In order to determine the matrix W , we optimize for

$$S_w = \sum_{k=1}^K \sum_{i, l_i=k} (\phi_i - \mu_k)(\phi_i - \mu_k)^T$$

$$S_b = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

where

$$\mu_k = \frac{1}{N_k} \sum_{i, l_i=k} \phi_i = \frac{1}{N_k} \Phi \mathbf{1}_k$$

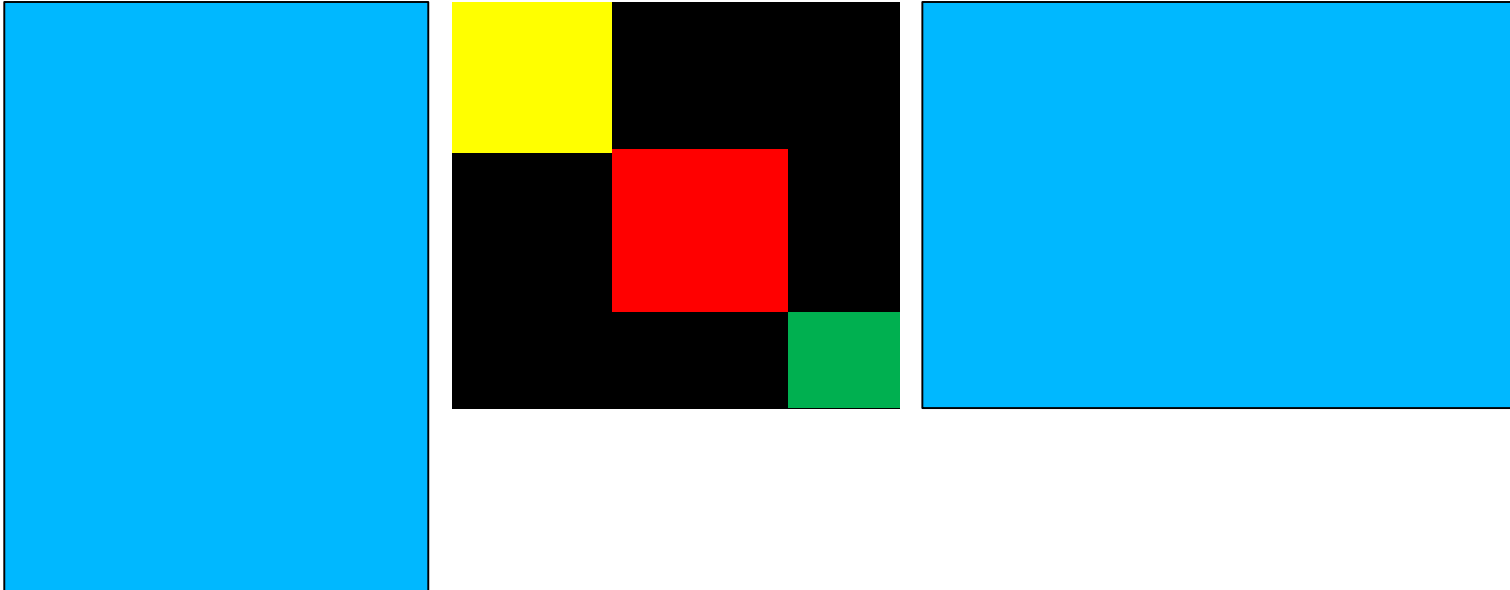
$$\mu = \frac{1}{N} \sum_{i=1}^N \phi_i = \frac{1}{N} \Phi \mathbf{1}$$

Kernel Discriminant Analysis

Substituting μ_k and μ in the scatter matrices, we get

$$\begin{aligned} S_w &= \sum_{k=1}^K \left(\Phi J_k - \frac{1}{N_k} \Phi \mathbf{1}_k \mathbf{1}_k^T \right) \left(\Phi J_k - \frac{1}{N_k} \Phi \mathbf{1}_k \mathbf{1}_k^T \right)^T \\ &= \sum_{k=1}^K \Phi \left(J_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right) \left(J_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right)^T \Phi^T \\ &= \Phi \left(\sum_{k=1}^K \left(J_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right) \left(J_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right)^T \right) \Phi^T \\ &= \Phi \mathbf{L}_w \Phi^T, \end{aligned}$$

Kernel Discriminant Analysis

$$\mathbf{S}_w = \Phi \mathbf{L}_w \Phi^T$$



The diagram illustrates the equation $\mathbf{S}_w = \Phi \mathbf{L}_w \Phi^T$. It shows three matrices arranged horizontally, separated by equals signs. The first matrix, Φ , is a large blue rectangle. The second matrix, \mathbf{L}_w , is a 4x4 block matrix with a black background. It contains a yellow block in the top-left corner, a red block in the center, and a green block in the bottom-right corner. The third matrix, Φ^T , is a large blue rectangle.

Kernel Discriminant Analysis

Substituting μ_k and μ in the scatter matrices, we get

$$\begin{aligned} S_b &= \sum_{k=1}^K N_k \left(\frac{1}{N_k} \Phi \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \Phi \mathbf{1} \mathbf{1}^T \right) \left(\frac{1}{N_k} \Phi \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \Phi \mathbf{1} \mathbf{1}^T \right)^T \\ &= \Phi \left(\sum_{k=1}^K N_k \left(\frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right) \left(\frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right)^T \right) \Phi^T \\ &= \Phi \mathbf{L}_b \Phi^T. \end{aligned}$$

Kernel Discriminant Analysis

$$\mathbf{S}_b = \Phi \mathbf{L}_b \Phi^T$$


The diagram illustrates the equation $\mathbf{S}_b = \Phi \mathbf{L}_b \Phi^T$. It shows three matrices arranged horizontally, separated by equals signs. The first matrix, Φ , is a large blue rectangle. The second matrix, \mathbf{L}_b , is a green square with black diagonal blocks. The third matrix, Φ^T , is another large blue rectangle.

Kernel Discriminant Analysis

Thus, the optimization criterion becomes

$$\mathcal{J}(\mathbf{W}) = \frac{\text{Tr}(\mathbf{W}^T (\Phi \mathbf{L}_b \Phi^T) \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\Phi \mathbf{L}_w \Phi^T) \mathbf{W})}$$

or using $\mathbf{W} = \Phi \mathbf{A}$

$$\mathcal{J}(\mathbf{A}) = \frac{\text{Tr}(\mathbf{A}^T (\mathbf{K} \mathbf{L}_b \mathbf{K}^T) \mathbf{A})}{\text{Tr}(\mathbf{A}^T (\mathbf{K} \mathbf{L}_w \mathbf{K}^T) \mathbf{A})} = \frac{\text{Tr}(\mathbf{A}^T \mathbf{S}_b^{(A)} \mathbf{A})}{\text{Tr}(\mathbf{A}^T \mathbf{S}_w^{(A)} \mathbf{A})}$$

Kernel Discriminant Analysis

Thus, the optimization criterion becomes

$$\mathcal{J}(\mathbf{W}) = \frac{\text{Tr}(\mathbf{W}^T (\Phi \mathbf{L}_b \Phi^T) \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\Phi \mathbf{L}_w \Phi^T) \mathbf{W})}$$

or using $\mathbf{W} = \Phi \mathbf{A}$

$$\mathcal{J}(\mathbf{A}) = \frac{\text{Tr}(\mathbf{A}^T (\mathbf{K} \mathbf{L}_b \mathbf{K}^T) \mathbf{A})}{\text{Tr}(\mathbf{A}^T (\mathbf{K} \mathbf{L}_w \mathbf{K}^T) \mathbf{A})} = \frac{\text{Tr}(\mathbf{A}^T \mathbf{S}_b^{(A)} \mathbf{A})}{\text{Tr}(\mathbf{A}^T \mathbf{S}_w^{(A)} \mathbf{A})}$$

which is solved by solving $\mathbf{S}_b^{(A)} \mathbf{a} = \lambda \mathbf{S}_w^{(A)} \mathbf{a}$

Kernel Discriminant Analysis

Demo SVM:

<http://cs.stanford.edu/people/karpathy/svmjs/demo/>

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

Demo MLP:

<http://cs.stanford.edu/~karpathy/svmjs/demo/demonn.html>