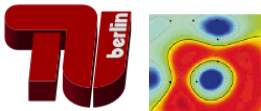


Lecture 11: Support Vector Machines

Machine Learning 1



Hyperplane $y = \text{sgn}(\mathbf{w} \cdot \Phi(x) + b)$ in \mathcal{F}

$$\begin{array}{ll} \min & \|\mathbf{w}\|^2 \\ \text{subject to} & y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] \geq 1 \quad \text{for } i = 1 \dots N \end{array}$$

(i.e. training data separated correctly, otherwise introduce slack variables).

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i \cdot ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1).$$

obtain unique α_i by QP (no local minima!): **dual problem**

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0, \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0,$$

$$\text{i.e.} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i).$$

Substitute both into L to get the **dual problem**



Hyperplane in \mathcal{F} with slack variables: SVM

$$\min \quad \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i^p$$

$$\text{subject to} \quad y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for } i = 1 \dots N$$

(introduce slack variables if training data **not** separated correctly)

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i \cdot ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1).$$

obtain unique α_i by QP (no local minima!): **dual problem**

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0, \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0,$$

$$\text{i.e.} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i).$$

Substitute both into L to get the **dual problem**



Dual Problem

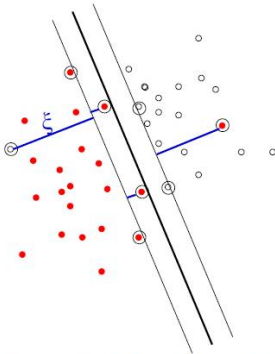
$$\begin{aligned} \text{maximize} \quad & W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad & C \geq \alpha_i \geq 0, \quad i = 1, \dots, N, \quad \text{and} \quad \sum_{i=1}^N \alpha_i y_i = 0. \end{aligned}$$

Note: solution determined by training examples (SVs) on /in the margin. Remark: duality gap.

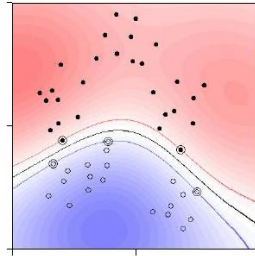
$$\begin{aligned} y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] &> 1 && \implies \alpha_i = 0 \longrightarrow \mathbf{x}_i \text{ irrelevant or} \\ y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] &= 1 && (\text{on /in margin}) \longrightarrow \mathbf{x}_i \text{ Support Vector} \end{aligned}$$



A Toy Example: $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2)$



linear SV with slack variables



nonlinear SVM, Domain: $[-1, 1]^2$

Kernel Trick

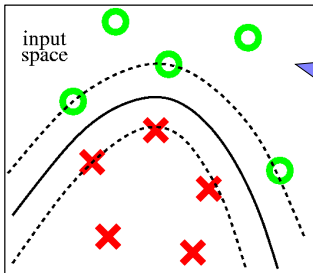
- Saddle Point: $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i)$.
- Hyperplane in \mathcal{F} : $y = \text{sgn}(\mathbf{w} \cdot \Phi(x) + b)$
- putting things together “kernel trick”

$$\begin{aligned} f(\mathbf{x}) &= \text{sgn}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b\right) \\ &= \text{sgn}\left(\sum_{i \in \#SVs} \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i) + b\right) \quad \text{sparse!} \end{aligned}$$

- trick: $k(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})$, i.e. **never use Φ : only k !!!**



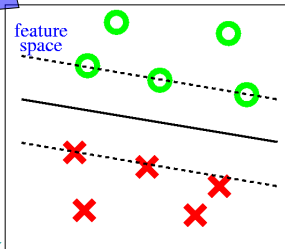
Support Vector Machines in a nutshell



good theory
non-linear decision by
implicitly **mapping** the data
into feature space by SV **kernel** function **K**

$$\Phi \text{ rsp. } K(x,y) = \Phi(x) \cdot \Phi(y)$$

Φ



Kernels ...

- kernels hold key to learning problem.
- **choosing kernels ...**
 - Mercer condition (ℓ_2 integrability & positivity)
 - kernel reflects prior (Smola, Schölkopf & Müller 98, Girosi 98)
 - approximating LOO bounds give good model selection results (Tsuda et al. 2001, Vapnik & Chapelle 2000)
- So: **engineer** an appropriate kernel from prior knowledge! (Jaakola and Haussler 1998, Watkins 2000, Zien et al 2000, recently a large body of interesting work)
- And: use **careful** model selection to find appropriate kernel parameters, i.e. chose appropriate degree of polynomial or bandwidth of Gaussian kernel



Digestion: Use of kernels

- **Question:** What makes kernel methods (e.g. SVM) perform well?
- **Answer:**
 - In the first place: a good idea/theory.
 - But also: **The kernel**
- Using kernels, we work explicitly in extremely high dimensional spaces (RKHS) with interesting features for themselves (depending on the kernel) [SSM et al. 98]
- Common choices: Gaussian kernel $\exp(-\|x-y\|^2/c)$ or polynomial kernel $(x \cdot y)^d$.
- Almost any linear algorithm can be transformed to feature space. [SSM et al. 98]
- With suitable regularization it outperforms its linear counterpart. [Mika et al. 02]
[Zien et al. 00, Tsuda et al. 02, Sonnenburg et al. 05]
- **The kernel can be adopted to specific tasks, e.g. using prior knowledge**



Kernels for graphs,
trees, strings etc.

Digestion

$$R[f] \leq R_{emp}[f] + \sqrt{\frac{d(\log \frac{2N}{d} + 1) - \log(\eta/4)}{N}}$$
$$K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})$$

