

Makeup class : Friday 31<sup>st</sup>

Select time on Canvas

Learning with kernels, SVMs with Gaussian kernels on MNIST data  
homework 1 (Sunday) (USPS)

## Goal

- Nonlinearly separable datasets
- how to extend SVMs & kernel SVMs

## Vanilla SVM

$$\min_{\omega, b} \frac{1}{2} \|\omega\|^2 \quad \text{s.t.} \quad y_i(\omega \cdot x_i + b) \geq 1 \quad \forall i \in [m]$$

$m$ : # of data pts

$d$ :  $\dim(\mathcal{X})$

$$\text{Soln: } \omega = \sum_{i=1}^m \lambda_i y_i x_i$$

$\lambda_i \neq 0$ ,  $x_i$ : support vectors

$\lambda = [\lambda_1, \dots, \lambda_m] \in \mathbb{R}^m$  dual variables

Lagrangian:

$$\mathcal{L}(\omega, b, \lambda) = \frac{\|\omega\|^2}{2} + \sum_{i=1}^m \lambda_i (1 - y_i(\omega \cdot x_i + b))$$

Dual problem: obtained by plugging in KKT conditions

$$\omega = \sum_{i=1}^m \lambda_i y_i x_i$$

$$\sum_{i=1}^m y_i \lambda_i = 0$$

$$\lambda_i = 0 \text{ or } y_i(\omega \cdot x_i + b) = 1$$

$$\text{Dual problem: } \max_{\lambda} \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \lambda_i \lambda_j y_i y_j x_i \cdot x_j$$

$$\lambda_i \geq 0 \quad i \in [m]$$

In dot product form;

Replace  $x_i \cdot x_j \rightarrow \Phi(x_i) \cdot \Phi(x_j)$

kernel SVM:

$$\max_{\lambda} \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \lambda_i \lambda_j y_i y_j \chi(x_i, x_j)$$

$$h_{\chi\text{-svm}}(x) = \text{sgn}\left(\sum_{i=1}^m \lambda_i y_i \chi(x_i, x) + b\right)$$

To solve for  $b$ :

for any  $i$  at which  $\lambda_i \neq 0$ ,

from KKT conditions,

$$y_i(\omega \cdot \Phi(x_i) + b) = 1$$

$$\omega \cdot \Phi(x_i) + b = y_i$$

$$b = y_i - \omega \cdot \Phi(x_i)$$

$$= y_i - \Phi(x_i) \cdot \sum_{j=1}^m \lambda_j y_j \Phi(x_j)$$

$$b = y_i - \sum_{j=1}^m \lambda_j y_j \chi(x_i, x_j)$$

## Kernelization of the SVM

$$x \rightarrow \Phi(x)$$

data features

$$\dim d \quad D \gg d$$

linear model on feature space

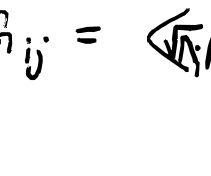
equivalent to

nonlinear model on data space

Recap: Cover's Theorem:

in higher dim space, data are more likely to be linearly separable.

MNIST dat set



$$28 \times 28 = 784$$

$$d = 784$$

intrinsic dimension:

$$O(10)$$

Linear classifiers not better than random guess!

## Recap

Gaussian kernel  $\left. \begin{array}{l} \text{PD} \\ \text{PD} \end{array} \right\} e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$  hyperparameter  $\sigma$

Polynomial kernel  $(x \cdot x' + c)^k$   $c, k$

Sigmoid kernel  $\tanh(a \cdot x \cdot x' + b)$   $a, b$

$b, a < 0$

not PD

PD kernel:  $\chi$  is symmetric

Gram matrix  $G_{ij} = \chi(x_i, x_j)$

is PSD. for all  $m \in \mathbb{N}$ .

Merz's Theorem

If

$$\chi(x, x') = \langle \Phi(x), \Phi(x') \rangle \text{ (Euclidean or } \ell^2 \text{ inner product)}$$

then  $\chi$  is PD

Converse also holds.

Proof: if  $\chi$  is s.t.  $\chi(x, x') = \langle \Phi(x), \Phi(x') \rangle$

then,

$$c^T G c \geq 0 \text{ for any } c \in \mathbb{R}^m.$$

$$\sum_{i,j=1}^m c_i c_j G_{ij} = \sum_{i,j=1}^m c_i c_j \Phi(x_i) \cdot \Phi(x_j)$$

$$= \left\| \sum_{i=1}^m c_i \Phi(x_i) \right\|^2 \geq 0.$$

Only if:

$G$  is SPSPD.

$$c^T G c \geq 0 \text{ for any } c \in \mathbb{R}^m.$$

$$G = P^T \Lambda P \quad P^T = P^{-1} \quad \Lambda_i \geq 0$$

$$G_{ij} = \langle \sqrt{\Lambda_i} p_i, \sqrt{\Lambda_j} p_j \rangle$$

$$\begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} \begin{bmatrix} \Lambda_1 & \\ & \Lambda_2 \end{bmatrix} \begin{bmatrix} p_{11} & p_{21} \\ p_{12} & p_{22} \end{bmatrix}$$

$$\begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} \begin{bmatrix} \Lambda_1 p_{11} & \Lambda_1 p_{21} \\ \Lambda_2 p_{12} & \Lambda_2 p_{22} \end{bmatrix} = \begin{bmatrix} \sqrt{\Lambda_1} p_{11}^2 \sqrt{\Lambda_1} & + \sqrt{\Lambda_2} \sqrt{\Lambda_1} p_{12}^2 \\ \sqrt{\Lambda_1} p_{11} \sqrt{\Lambda_1} p_{21} + \sqrt{\Lambda_2} \sqrt{\Lambda_2} p_{22} p_{12} \end{bmatrix}$$

$$P_{21} \Lambda_1 p_{11} + P_{22} \Lambda_2 p_{12}$$

$$\sqrt{\Lambda_1} p_{11} \sqrt{\Lambda_1} p_{21} + \sqrt{\Lambda_2} \sqrt{\Lambda_2} p_{22} p_{12}$$

$$\Phi(x_i) = \sqrt{\Lambda_i} p_i \rightarrow \text{ith row of matrix } P \text{ whose columns are eigenvectors of } \Phi(\Phi)$$

(Remember analog in  $\infty$  dims)

$$\Phi(x_i) \cdot \Phi(x_j) = G_{ij}$$

features

$\Phi$  are defined only at data points but they may not be "linearly" extensible to  $\mathcal{X}$ .

Even if  $\chi$  is not PD but  $G$  is SPSPD, we may still use features (extracted above) to "kernelize" algorithms in dot product form

## Nonlinearly separable data

$$\nexists \omega, b \text{ s.t. } y_i(\omega \cdot x_i + b) > 0 \text{ at all data pts } x_i, y_i$$

$$y_i(\omega \cdot \Phi(x_i) + b) > 0$$

## Soft-margin SVM

$$\min_{\omega, \xi, b} \frac{\|\omega\|^2}{2} + C \|\xi\|^p$$

$$\text{s.t. } y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

(Vanilla SVM / hard-margin SVM: margin  $\frac{1}{\|\omega\|}$ )



$$\min_{\omega, \xi, b} \frac{\|\omega\|^2}{2} + C \|\xi\|^p \quad \text{obj. of soft margin classification}$$

if  $C$  is large, we penalize "outliers"

if  $C$  is small, we focus on max. margin

$$P=1 \text{ or } P=2$$

$$\min_{\omega, \xi, b} \frac{\|\omega\|^2}{2} + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

$$\mathcal{L}(\omega, b, \xi, \lambda, \gamma) = \frac{\|\omega\|^2}{2} + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \lambda_i (y_i(\omega \cdot x_i + b) - 1 + \xi_i) - \sum_{i=1}^m \gamma_i \xi_i$$

KKT

$$\omega = \sum_{i=1}^m \lambda_i y_i x_i$$

$$\sum_{i=1}^m \lambda_i y_i = 0$$

$$C = \lambda_i + \gamma_i$$

$$y_i(\omega \cdot x_i + b) = 1 - \xi_i \text{ or } \lambda_i = 0$$

$$\gamma_i = 0 \text{ or } \xi_i = 0$$

Ex: kernel soft margin SVM dual problem

Support vectors  $x_i$  for which  $\lambda_i \neq 0$ .

$$y_i(x \cdot \omega + b) = 1 - \xi_i$$

Two types of support vectors:

$\xi_i = 0$  Location of  $x_i$  marginal hyperplane

$\xi_i \neq 0$  outlier

$$\lambda_i = C$$

