

SVM: Support Vector Machine (classification)

Linear

Constructing a projection

- Datapoint  $x$
- Projection vector  $w$  (Normal to plane)
- projection of  $x$  onto  $w$   
 $w^T \cdot x = \|w\| \|x\| \cos(\theta)$
- $\Rightarrow$  the sign allows to separate points on either side of the plane  
 $w^T \cdot x > 0 \rightarrow$  on the left handside of plane  
 $w^T \cdot x < 0 \rightarrow$  on the right handside of plane

Linear classifiers

- Class label  $y = \{-1; 1\}$
- Separating hyperplane is comprised of normal and intercept  
 $w$ : plane normal  
 $b$ : intercept
- Decision function  $y = f(x; w, b) = \text{sgn}(w^T x + b)$

Classifier margin

- definition: the width of the boundary
- best measure of SVM: the maximum margin
- Support vectors
  - datapoints closest to the boundary
  - define the margin

Computing the Distance to the Separating Hyperplane

- The margins on either side of the hyperplane satisfy  $\langle w, x \rangle + b = \pm 1$ 
  - positive  $\{x : \langle w, x \rangle + b = +1\}$
  - negative  $\{x : \langle w, x \rangle + b = -1\}$
- Distance to place = projection
  - distance projection with unitary vector  $\frac{\langle w, (x' - x) \rangle}{\|w\|} \frac{w}{\|w\|}$
  - datapoint porjection on the plane  $\langle w, x' \rangle + b = 0 \Rightarrow \langle w, x' \rangle = -b$
  - distance to place  $\frac{|\langle w, x \rangle + b|}{\|w\|}$

Computing the margin

- Distance of each points on either side of the margin  $\Rightarrow$  margin between two classes
  - $\|x^1 - x^2\| = \frac{|\langle w, x^1 \rangle + b|}{\|w\|} = \frac{1}{\|w\|}$
  - $\|x^2 - x^2\| = \frac{|\langle w, x^2 \rangle + b|}{\|w\|} = \frac{1}{\|w\|}$ $\Rightarrow \|x^1 - x^2\| = \frac{2}{\|w\|}$

Objective function (Primal problem)

- minimize the convex form  $\frac{\|w\|^2}{2}$  (Maximize  $\Rightarrow$  Minimize)
- constraints  $y^i (\langle w, x^i \rangle + b) \geq 1, i = 1, 2, \dots, M$
- Misc.
  - $(N+1)$  parameters  $w \in \mathbb{R}^N, b \in \mathbb{R}$  — dimension of data
  - $M$  constraints — number of data

Solving the constrained optimization

- Lagrange method
  - introduction  $M$  Lagrange multipliers  $\alpha$
  - MaxMin problem — maximizing over  $\alpha$  and minimizing over  $w$  and  $b$
  - Requesting the gradient
- Solution
  - Normal vector  $\frac{\partial L(w, b, \alpha)}{\partial w} = 0 \Leftrightarrow w = \sum_{i=1}^M \alpha_i y^i x^i$  (global solution but not identity, different linear combinations lead to different solutions)
  - Intercept  $\frac{\partial L(w, b, \alpha)}{\partial b} = 0 \Leftrightarrow \sum_{i=1}^M \alpha_i y^i = 0$

Dual optimization problem

- Taking the definition of  $w$  and plugging it back to the Lagrangian
- Solved through Sequential Minimal Optimization algorithm (SMO)
- Objective function  $\max_{\alpha} W(\alpha_i) = \sum_{i=1}^M \alpha_i - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j y_i y_j x_i^T x_j$
- Constraints subject to:  $\alpha_i \geq 0$  and  $\sum_{i=1}^M \alpha_i y_i = 0$

The KKT Conditions

- Ensure that the primal and dual optimization problems have the same optimal solutions
  - Support vectors  $\alpha_i > 0$
  - Ignored points  $\alpha_i = 0$
- Decision function in terms of the support vectors  $f(x) = \text{sgn} \left( \sum_{i=1}^M \alpha_i y^i \langle x, x^i \rangle + b \right)$
- Use KKT to compute  $b \left( y^j \left( \left\langle \sum_{i=1}^M \alpha_i y^i x^i, x^j \right\rangle + b \right) - 1 \right) = 0$

Non-separable datasets

- Introduce slack variables  $\xi_i \geq 0$ 
  - correct classification  $\{0\}$
  - correct classification inside margin  $(0, 1)$
  - missclassification  $[1, \infty)$
- Relax the constraints  $y_i (\langle w, x \rangle + b) \geq 1 - \xi_i$
- Modify the optimization function  $\min_{w, b, \xi} \left( \frac{1}{2} \|w\|^2 + \frac{C}{M} \sum_{j=1}^M \xi_j \right)$  (Tradeoff: The higher  $C$  is, the better fitting the model is  $\Rightarrow$  lead to over-fitting possibly,  $\frac{C}{M}$ )

Non-linear

Using polar coordinates  $\Rightarrow$  Data become linearly separable — e.g. two groups of points are randomly distributed with different radii

data  $X$  — feature space  $H$  — nonlinear map  $\phi$

Map original input space (data  $X$ ) into high-dimension feature space through the nonlinear map — Idea: Making the problem linear

The dimension of mapped feature space may be greater than original one

Determining  $\phi$  is difficult and even impossible — Finding the transformation

Sufficient to compute distance in feature space

Kernel function  $k : X \times X \rightarrow \mathbb{R}$   
 $k(x^i, x^j) \rightarrow \langle \phi(x^i), \phi(x^j) \rangle$  — Kernel method

RBF (Gaussian) kernel  $k(x, x^i) = \exp(-\frac{\|x - x^i\|^2}{2\sigma^2})$

Polynomial kernel  $k(x, x^i) = (\langle x, x^i \rangle + c)^p$  — Different kernels

Hyperparameters that need to be determined by hand

Key components

- Hyperplane
- Support vectors
- Margin — Region of equal distance on both side
- Color gradient — Distance to hyperplane

Effect of hyperparameters on performance

- $C$  that determines the costs associated to incorrectly classifying datapoints
- Higher  $C \Rightarrow$  Higher accuracy (possible over-fitting)

Multi-class classification

- Construct a set of  $K$  binary classifiers
- Sufficient to use  $K-1$  classifier for  $K$  classes
- Computing the  $K$ -th classifier may provide tighter bounds

Misc.

No confidence in prediction

- Does not entail a notion of confidence! (no a notion of likelihood)
- Predict by default sign of  $b$  lead to large amount of false positives
- Doing crossvalidation would not prevent the effect as it is close to training datapoints.
- Verify the sign of  $b$  for specific class
- Improvements
  - Run crossvalidation by generating the datapoints that never seen
  - Compare the distribution of train/ test

Summary

- Pros
  - Compared to KNN — Build a model & boundary is a specific function
  - Compared to GMM — Guarantee to find global optimum
- Cons
  - Computing costs at testing  $O(M)$ ,  $M$  is number of datapoints
  - Require to choose two more parameters (Compared to 1 for KNN and GMM)
  - Weight of slack variable  $C$
  - Kernel width  $\sigma$