

# STA141C: Big Data & High Performance Statistical Computing

## Lecture 9: Classification

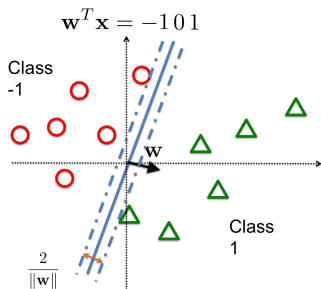
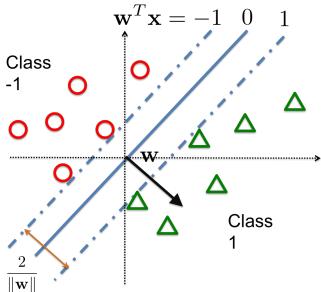
Cho-Jui Hsieh  
UC Davis

May 18, 2017

# Linear SVM

# Support Vector Machines

- SVM is a widely used classifier.
- Given:
  - Training data points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .
  - Each  $\mathbf{x}_i \in \mathbb{R}^d$  is a feature vector:
  - Consider a simple case with two classes:  $y_i \in \{+1, -1\}$ .
- Goal: Find a hyperplane to separate these two classes of data:  
if  $y_i = 1$ ,  $\mathbf{w}^T \mathbf{x}_i \geq 1$ ; if  $y_i = -1$ ,  $\mathbf{w}^T \mathbf{x}_i \leq -1$ .



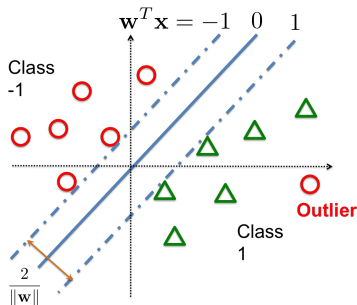
# Support Vector Machines (hard constraints)

- Given training data  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  with labels  $y_i \in \{+1, -1\}$ .
- SVM primal problem (with hard constraints):

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i) \geq 1, i = 1, \dots, n,$$

- What if there are outliers?

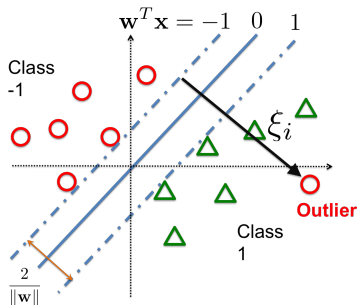


# Support Vector Machines

- Given training data  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  with labels  $y_i \in \{+1, -1\}$ .
- SVM primal problem:

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i) \geq 1 - \xi_i, i = 1, \dots, n,$$
$$\xi_i \geq 0$$



# Support Vector Machines

- SVM primal problem can be written as

$$\min_{\mathbf{w}} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{L2 regularization}} + \sum_{i=1}^n \underbrace{\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)}_{\text{hinge loss}}$$

- Non-differentiable when  $y_i \mathbf{w}^T \mathbf{x}_i = 1$  for some  $i$

# General Empirical Risk Minimization

- Regularized ERM:

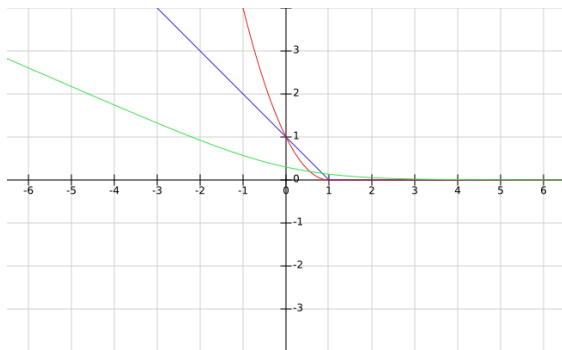
$$\min_{\mathbf{w} \in \mathbb{R}^d} P(\mathbf{w}) := \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) + R(\mathbf{w})$$

- $\ell_i(\cdot)$ : loss function
- $R(\mathbf{w})$ : regularization
- Dual problem may have a different form (?)

# Examples

- Loss functions:

- Regression:  $l_i(\mathbf{x}_i) = (\mathbf{x}_i - y_i)^2$
- SVM (hinge loss):  $l_i(\mathbf{x}_i) = \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)$
- Square hinge loss:  $l_i(\mathbf{x}_i) = \max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)^2$
- Logistic regression:  $l_i(\mathbf{x}_i) = \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$

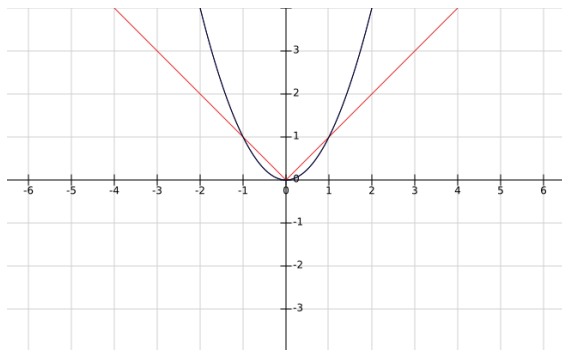




# Examples

- Regularizations:

- L2-regularization:  $\|\mathbf{w}\|_2^2$ : small but dense solution
- L1-regularization:  $\|\mathbf{w}\|_1$ : sparse solution
- Nuclear norm:  $\|W\|_*$ : low-rank solution



# LIBLINEAR

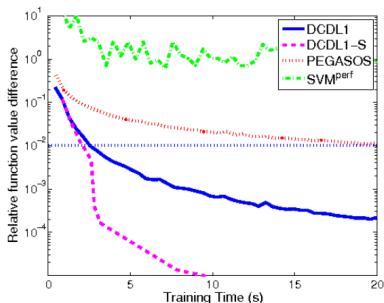
- Implemented in LIBLINEAR:

`https://www.csie.ntu.edu.tw/~cjlin/liblinear/`

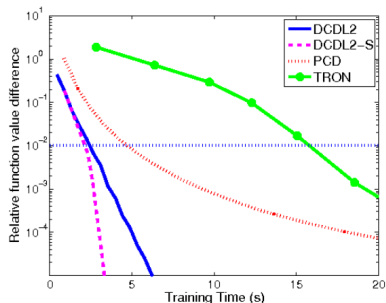
- Other functionalities:

- Logistic regression (L1 or L2 regularization)
- Multi-class SVM
- Support vector regression
- Cross-validation

- RCV1: 677,399 training samples; 47,236 features; 49,556,258 nonzeros in the whole dataset.



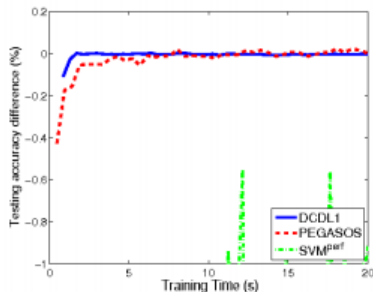
(e) L1-SVM: rcv1



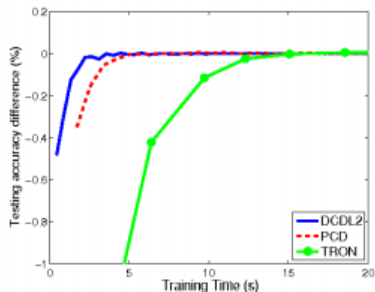
(f) L2-SVM: rcv1

Time vs primal objective function value

- RCV1: 677,399 training samples; 47,236 features; 49,556,258 nonzeros in the whole dataset.



(e) L1-SVM: rcv1



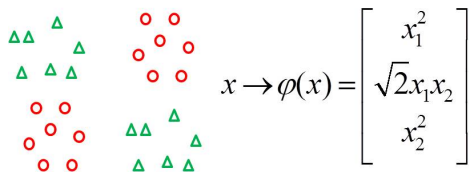
(f) L2-SVM: rcv1

Time vs prediction accuracy

# Kernel SVM

# Non-linearly separable problems

- What if the data is not linearly separable?



**Solution:** map data  $x_i$  to higher dimensional (maybe infinite) feature space  $\varphi(x_i)$ , where they are linearly separable.

# Support Vector Machines (SVM)

- SVM primal problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \varphi(\mathbf{x}_i)) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, n, \end{aligned}$$

- The dual problem for SVM:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha, \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad \text{for } i = 1, \dots, n, \end{aligned}$$

where  $Q_{ij} = y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$  and  $\mathbf{e} = [1, \dots, 1]^T$ .

- Kernel trick: define  $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ .
- At optimum:  $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \varphi(\mathbf{x}_i)$ ,

## Various types of kernels

- Gaussian kernel:  $K(\mathbf{x}_i, \mathbf{y}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2}$ ;
- Polynomial kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + c)^d$ .
- Hard to solve: need to solve  $n$ -by- $n$  quadratic minimization problem,  $\geq O(n^2)$  time.
- LIBSVM: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- For linear SVM, use LIBLINEAR instead of LIBSVM.



- Linear SVM: `sklearn.svm.LinearSVC`
- Logistic Regression: `sklearn.linear_model.LogisticRegression`
- Kernel SVM: `sklearn.svm.SVC`
- ...

- Linear SVM: `sklearn.svm.LinearSVC`
- Logistic Regression: `sklearn.linear_model.LogisticRegression`
- Kernel SVM: `sklearn.svm.SVC`
- ...
- Practice in homework.

# Coming up

- Classification

Questions?