#### Mollabashi house Isfahan





Introduction to Optimization for Machine Learning

Hossein Moosaei, PhD

Postdoctoral Researcher Department of Applied Mathematics School of Computer Science Charles University

4 November 2019

## Optimization

#### Finding the minimizer of a function subject to constraints:

 $egin{aligned} & \min_x \ f_0(x) \ & ext{s.t.} \ f_i(x) \leq 0, \ i = \{1, \dots, k\} \ & h_j(x) = 0, \ j = \{1, \dots, l\} \end{aligned}$ 

## Some different types Of optimization problems

#### **Optimization Taxonomy**



### Applications of Optimization



### Karush-Kuhn-Tucker Optimality Conditions

## **Optimality** Criteria

• <u>Big question</u>: How do we know that we have found the "optimum" for min f(x)?

Answer: Test the solution for the "necessary and sufficient conditions"

## Optimality Conditio<del>ns – Unconstrain</del>ed Case

• Let x\* be the point that we think is the minimum for f(x) Necessary condition (for optimality):

 $\nabla f(x^*) = 0$ 

- A point that satisfies the necessary condition is a stationary point It can be a minimum, maximum, or saddle point
- How do we know that we have a minimum?
- <u>Answer</u>: <u>Sufficiency Condition</u>:

The sufficient conditions for x\* to be a strict local minimum are:

 $\nabla f(x^*) = 0$  $\nabla^2 f(x^*)$  is positive definite

## Constrained Case – KKT Conditions

- To proof a claim of optimality in constrained minimization (or maximization), we have to check the found point with respect to the (Karesh) Kuhn Tucker conditions.
- Kuhn and Tucker extended the Lagrangian theory to include the general classical single-objective nonlinear programming problem:

 $\begin{array}{ll} \mbox{minimize} & f(x) \\ \mbox{Subject to} & g_j(x) \geq 0 \mbox{ for } j = 1, \, 2, \, \dots, \, J \\ & h_k(\textbf{x}) = 0 & \mbox{for } k = 1, \, 2, \, \dots, \, K \\ & x = (x_1, \, x_2, \, \dots, \, x_N) \end{array}$ 

## Necessary KKT Conditions

For the problem: Min f(x)s.t.  $g(x) \le 0$ (n variables, m constraints)

The necessary conditions are:  $\nabla f(x) + \Sigma \mu_i \nabla g_i(x) = 0$  (optimality)  $g_i(x) \le 0$  for i = 1, 2, ..., m (feasibility)  $\mu_i g_i(x) = 0$  for i = 1, 2, ..., m (complementary slackness condition)  $\mu_i \ge 0$  for i = 1, 2, ..., m (non-negativity)

Note that the first condition gives n equations.

# Necessary KKT Conditions (General Case)

 For general case (n variables, M Inequalities, L equalities): Min f(x)

s.t.

 $\begin{array}{ll} g_i(x) \leq 0 \mbox{for } i = 1, \, 2, \, ..., \, M \\ h_j(x) = 0 & \mbox{for } J = 1, \, 2, \, ..., \, L \end{array}$ 

- In all this, the assumption is that  $\nabla g_j(x^*)$  for j belonging to active constraints and  $\nabla h_k(x^*)$  for k = 1, ..., K are linearly independent
- The necessary conditions are:  $\nabla f(x) + \Sigma \mu_i \nabla g_i(x) + \Sigma \lambda_j \nabla h_j(x) = 0$  (optimality)  $g_i(x) \le 0$  for i = 1, 2, ..., M (feasibility)  $h_j(x) = 0$  for j = 1, 2, ..., L (feasibility)  $\mu_i g_i(x) = 0$  for i = 1, 2, ..., M (complementary slackness condition)  $\mu_i \ge 0$  for i = 1, 2, ..., M (non-negativity) (Note:  $\lambda_i$  is unrestricted in sign)

## Restating the Optimization Problem

• Kuhn Tucker Optimization Problem: Find vectors  $x_{(Nx1)}$ ,  $\mu_{(1xM)}$  and  $\lambda_{(1xK)}$  that satisfy:  $\nabla f(x) + \Sigma \mu_i \nabla g_i(x) + \Sigma \lambda_j \nabla h_j(x) = 0$  (optimality)  $g_i(x) \le 0$  for i = 1, 2, ..., M (feasibility)  $h_j(x) = 0$  for j = 1, 2, ..., L (feasibility)  $\mu_i g_i(x) = 0$  for i = 1, 2, ..., M (complementary slackness condition)  $\mu_i \ge 0$  for i = 1, 2, ..., M (non-negativity)

- > If x\* is an optimal solution to NLP, then there exists a  $(\mu^*, \lambda^*)$  such that  $(x^*, \mu^*, \lambda^*)$  solves the Kuhn–Tucker problem.
- Above equations not only give the necessary conditions for optimality, but also provide a way of finding the optimal point.

## Limitations

 <u>Necessity theorem</u> helps identify points that are not optimal. A point is not optimal if it does not satisfy the Kuhn–Tucker conditions.

- On the other hand, not all points that satisfy the Kuhn-Tucker conditions are optimal points.
- The Kuhn–Tucker <u>sufficiency theorem</u> gives conditions under which a point becomes an optimal solution to a single-objective NLP.

## Sufficiency Condition

- Sufficient conditions that a point x\* is a strict local minimum of the NLP problem, where f, g<sub>j</sub>, and h<sub>k</sub> are twice differentiable functions are that
  - 1) The necessary KKT conditions are met.
  - 2) The Hessian matrix  $\nabla^2 L(x^*) = \nabla^2 f(x^*) + \Sigma \mu_i \nabla^2 g_i(x^*) + \Sigma \lambda_j \nabla^2 h_j(x^*)$  is positive definite on a subspace of  $\mathbb{R}^n$  as defined by the condition:
    - $y^T \nabla^2 L(x^*) y \ge 0$  is met for every vector  $y_{(1xN)}$  satisfying:
      - $\nabla g_j(\mathbf{x^*})y < 0$  for j belonging to  $I_1 = \{j \mid g_j(\mathbf{x^*}) = 0, u_j^* > 0\}$  (active constraints)

 $\nabla h_k(\mathbf{x^*})y = 0 \text{ for } k = 1, ..., K$ 

y ≠ 0

## KKT Sufficiency Theorem (Special

## Case)

- Consider the classical single objective NLP problem.
  - $\begin{array}{ll} \mbox{minimize} & f(x) \\ \mbox{Subject to } g_j(x) \leq 0 & \mbox{for } j = 1, \, 2, \, ..., \, J \\ & h_k({\bm x}) = 0 & \mbox{for } k = 1, \, 2, \, ..., \, K \end{array}$
- Let the objective function f(x) be convex, the inequality constraints g<sub>j</sub>(x) be all convex functions for j = 1, ..., J, and the equality constraints h<sub>k</sub>(x) for k = 1, ..., K be linear.
- If this is true, then the necessary KKT conditions are also sufficient.
- Therefore, in this case, if there exists a solution x\* that satisfies the KKT necessary conditions, then x\* is an optimal solution to the NLP problem.
- In fact, it is a <u>global</u> optimum.

## **Dual Problem**

#### Generalized Lagrangian Function

 Consider the general (primal) optimization problem

minimize f(w)subject to  $g_i(w) \le 0, i = 1, \dots, k$  $h_i(w) = 0, j = 1, \dots, m$ 

where the functions f,  $g_i$ ,  $i = 1, \dots, k$ , and  $h_i$ ,  $i = 1, \dots, m$ , are defined on a domain  $\Omega$ . The generalized Lagrangian was defined as

$$L(w,\alpha,\beta) = f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{j=1}^{m} \beta_j h_j(w)$$
$$= f(w) + \alpha^T g(w) + \beta^T h(w)$$

#### Dual Problem and Strong Duality Theorem

 Given the primal optimization problem, the dual problem of it was defined as

> maximize  $\theta(\alpha,\beta) = \inf_{w \in \Omega} L(w,\alpha,\beta)$ subject to  $\alpha > 0$

• Strong Duality Theorem: Given the primal optimization problem, where the domain  $\Omega$  is convex and the constraints  $g_i$  and  $h_i$  are affine functions. Then the optimum of the primal problem occurs at the same values as the optimum of the dual problem .

## Machine Learning



### Unsupervised Learning

## What is Clustering?

Also called unsupervised learning, sometimes called classification by statisticians and sorting by psychologists and segmentation by people in marketing

Organizing data into classes such that there is •

- high intra-class similarity •
- low inter-class similarity •

Finding the class labels and the number of classes directly from • the data (in contrast to classification).



# What is a natural grouping among these objects?



# What is a natural grouping among these objects?

## Clustering is subjective





Simpson's Family School Employees





Females

Males

## A data set with clear cluster structure



 How would you design an algorithm for finding the three clusters in this case?

## **Supervised Learning**



## General Approach for Building Classification Model





## Supervised vs. Unsupervised Learning

- Supervised learning (classification)
  - Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
  - New data is classified based on the training set
- Unsupervised learning (clustering)
  - The class labels of training data is unknown
  - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

## **Classification Techniques**

#### • Base Classifiers

- Decision Tree based Methods
- Rule-based Methods
- Nearest-neighbor
- Neural Networks
- Deep Learning
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines

#### Ensemble Classifiers

• Boosting, Bagging, Random Forests

### Support Vector Machine (SVM)
## What is a good Decision Boundary?

• Consider a two-class, linearly separable classification problem. Construct the hyperplane  $w^{T}x+b=0, x \in R^{n}$ 

to make

$$w^{T} x_{i} + b > 0,$$
 for  $y_{i} = +1$   
 $w^{T} x_{i} + b < 0,$  for  $y_{i} = -1$ 

 Many decision boundaries! Are all decision boundaries equally good?



### Examples of Bad Decision Boundaries





## Optimal separating hyperplane

### • The optimal separating hyperplane



# • For the hyperplane, it can be proved that the margin *m* is

$$m = \frac{1}{\left\| w \right\|^2}$$

Hence, maximizing margin is equivalent to minimizing the square of the norm of  $\mathcal{W}$ .

### Finding the optimal decision boundary

- Let  $\{x_1, ..., x_n\}$  be our data set and let  $y_i \in \{1, -1\}$ be the class label of  $x_i$
- The optimal decision boundary should classify all points correctly  $\Rightarrow y_i(w^T x_i + b) \ge 1, \Box i$
- The decision boundary can be found by solving the following constrained optimization problem

minimize 
$$\frac{1}{2} \|w\|^2$$
  
subject to  $y_i (w^T x_i + b) \ge 1 \quad \forall i$ 

# Lagrangian of the optimization problem *minimize* $\frac{1}{2} \|w\|^2$ *subject to* $y_i(w^T x_i + b) \ge 1 \quad \forall i$

• The Lagrangian is

$$L = \frac{1}{2} w^{T} w + \sum_{i=1}^{n} \alpha_{i} (1 - y_{i} (w^{T} x_{i} + b))$$

• Setting the gradient of L w.r.t. w and be to zero, we have

$$w + \sum_{i=1}^{n} \alpha_i (-y_i) x_i = 0 \implies w = \sum_{i=1}^{n} \alpha_i y_i x_i$$
$$\sum_{i=1}^{n} \alpha_i y_i = 0$$

## The Dual Problem



• Note that  $\sum_{i=1}^{n} \alpha_{i} y_{i} = 0$ , and the data points appear in terms of their inner product; this is a quadratic function of  $\alpha_{i}$  only.

## The Dual Problem

• The dual problem is therefore:

maxmize 
$$W(\alpha) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j$$
  
subject to  $\alpha_i \ge 0$ ,  $\sum_{i=1}^{n} \alpha_i y_i = 0$ 

**The Dual Problem** minimize  $W(\alpha) = \frac{1}{2} \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{n} \alpha_i$ subject to  $\alpha_i \ge 0$ ,  $\sum_{i=1}^{n} \alpha_i y_i = 0$ • This is a quadratic programming (QP) problem, and therefore a global minimum of  $\alpha_i$  can always be found

• *w* can be recovered by  $w = \sum_{i=1}^{n} \alpha_{i} y_{i} x_{i}$ , and  $b = y_{k} - \sum_{i=1}^{n} \alpha_{i} y_{i} x_{i}^{T} x_{k}$  for any  $\alpha_{k} > 0$ 

so the decision function can be written

$$f(x) = sign\left(\sum_{i=1}^{n} \alpha_{i} y_{i} x_{i}^{T} x + b\right)$$

## The use of slack variables

• We allow "errors"  $\xi_i$  in classification for noisy data



## Soft Margin Hyperplane

The use of slack variables ξ<sub>i</sub> enable the soft margin classifier

$$\begin{cases} w^{T} x_{i} + b \ge 1 - \xi_{i} & y_{i} = 1 \\ w^{T} x_{i} + b \le -1 + \xi_{i} & y_{i} = -1 \\ \xi_{i} \ge 0 & \forall i \end{cases}$$

ξ<sub>i</sub> are "slack variables" in optimization
Note that ξ<sub>i</sub>=0 if there is no error for x<sub>i</sub>

• The objective function  $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$ *C* : tradeoff parameter between error and margin

### • The primal optimization problem becomes

minimize 
$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$
  
subject to  $y_i (w^T x_i + b) \ge 1 - \xi_i, \quad \xi_i \ge 0$ 

## **Dual Soft-Margin Optimization Problem**

• The dual of this new constrained optimization problem is maxmize  $W(\alpha) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j$ 

subject to 
$$C \ge \alpha_i \ge 0$$
,  $\sum_{i=1}^n \alpha_i y_i = 0$ 

• *w* can be recovered as  $w = \sum_{i=1}^{n} \alpha_i y_i x_i$ 

 This is very similar to the optimization problem in the hard-margin case, except that there is an upper bound C on α<sub>i</sub> now.

• Once again, a QP solver can be used to find  $\alpha_i$ 

## Nonlinear separable problems



## Non-linear SVMs: Feature spaces



### **Proximal Support Vector Machine**



The algorithm finds two non-parallel hyperplanes one for each class, each hyperplane should be as close as possible to one class and as far as possible from the other class.

 $\min \frac{\|AW^1 + b^1\|}{\|BW^1 + b^1\|}$ 

 $\min \frac{\|BW^2 + b^2\|}{\|AW^2 + b^2\|}$ 

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE,

### Twin Support Vector Machines for Pattern Classification

#### Jayadeva, *Senior Member*, *IEEE*, R. Khemchandani, *Student Member*, *IEEE*, and Suresh Chandra

Abstract—We propose Twin SVM, a binary SVM classifier that determines two nonparallel planes by solving two related SVM-type problems, each of which is smaller than in a conventional SVM. The Twin SVM formulation is in the spirit of proximal SVMs via generalized eigenvalues. On several benchmark data sets, Twin SVM is not only fast, but shows good generalization. Twin SVM is also useful for automatically discovering two-dimensional projections of the data.



# Why TWSVM?

This quadratic programming problem (QPP) is expensive to solve for large dimensions because all data points appear in the constraints.

## How does it works?

Instead of solving one large QPP, TWSVM solve two smaller OPP each of them has the formulation of standard SVM except that not all data patterns appear in the constraint at the same time. The algorithm finds two non-parallel hyperplanes one for each class, each hyperplane should be as close as possible to one class and as far as possible from the other class.



## **Linear Classifier**

TWSVM is obtained by solving the following pair of QPPs:

$$\begin{array}{ll} (TWSVM1) & \underset{w^{(1)}, b^{(1)}, q}{Min} & \frac{1}{2} (Aw^{(1)} + e_1 b^{(1)})^T (Aw^{(1)} + e_1 b^{(1)}) + c_1 e_2^T q \\ & subject \ to & - (Bw^{(1)} + e_2 b^{(1)}) + q \ge e_2, \ q \ge 0, \end{array}$$

$$\begin{array}{ll} (TWSVM2) & \underset{w^{(2)}, b^{(2)}, q}{Min} & \frac{1}{2} (Bw^{(2)} + e_2 b^{(2)})^T (Bw^{(2)} + e_2 b^{(2)}) + c_2 e_1^T q \\ & subject \ to & (Aw^{(2)} + e_1 b^{(2)}) + q \ge e_1, \quad q \ge 0, \end{array}$$

The first term of the objective function represents the sum of square distance from the hyperplane to each pattern of one class, therefore minimizing it keeps the hyperplane close to the patterns of one class.

The constraints require the hyper plane to be far from the other class patterns at least with distance 1.

The second term of the objective function minimize the sum of error variables to minimize miss classification of patterns belongs to other class. The Wolfe dual can be obtain as follows

$$\max_{\alpha} e_2^T \alpha - \frac{1}{2} \alpha^T G(^{H^T H) - 1} G^T \alpha, \qquad G = \begin{bmatrix} B & e_2 \end{bmatrix} \text{ and } H = \begin{bmatrix} A & e_1 \end{bmatrix}$$
  
subject to  $0 \le \alpha \le c_1$ 

 $u = -(H^T H)^{-1} G^T \alpha$  where  $u = [w_1^T, b_1]^T$ .

$$\max_{\alpha} e_1^T \gamma - \frac{1}{2} \gamma^T P(Q^T Q)^{-1} P^T \gamma, \qquad P = [A \ e_1] \text{ and } Q = [B \ e_2]$$
  
subject to  $0 \le \gamma \le c_2$ 

 $v = (Q^T Q)^{-1} P^T \gamma$  where  $v = [w_2^T, b_2]^T$ 

The first QPP TWSVM can be modified as follow:

$$\min_{\substack{w_1, b_1, q_1}} \|Aw_1 + e_1b_1\|^2 + c_1q^Tq,$$
  
subject to  $-(Bw_1 + e_2b_1) + q \ge e_2,$   
 $q \ge 0.$ 

We combine constraint together, then we have

$$q = (e_2 - Bw_1 - e_2b_1)_+$$

Then the above problem change to following unconstrained problem:

$$\min_{w_1,b_1,q_1} \|Aw_1 + e_1b_1\|^2 + c_1\|(e_2 - Bw_1 - e_2b_1)_+\|^2.$$

Similarly, the second QPP TWSVM can be modified as follow:

$$\min_{w_2,b_2,q_2} \|Bw_2 + e_2b_2\|^2 + c_2\|(e_1 - Aw_2 - e_1b_2)_+\|^2.$$

The above problems are piecewise, quadratic, convex, and once differentiable. The generalized Newton method can be used for solving them. Algorithm : Generalized Newton Method with the Stepsize Armijo Rule

Choose any vector  $p_0$  and  $\epsilon > 0$ , i = 0;

while  $\|\nabla g(p_i)\|_{\infty} \ge \epsilon$ 

Choose  $\alpha_i = max\{s, s\delta, s\delta^2, ...\}$  such that

$$g(p_i) - g(p_i + \alpha_i d_i) \ge -\delta \mu \nabla g(p_i)^T d_i,$$

where  $d_i = -\partial^2 g(p_i)^{-1} \nabla g(p_i)$ , s > 0 is a constant,  $\delta \in (0, 1)$  and  $\mu \in (0, 1)$ .

 $p_{i+1} = p_i + \alpha_i d_i$ i = i + 1;

#### **Numerical Experiments**

Data set	Twin SVM	New Method	
ionosphere	0.8346+-0.0617	0.92024e-001+- 3.9924e-002	
WPBC	0.6511+-0.2512	0.8792+-0.0757	
WDBC	0.5778+-0.1128	0.9526+-0.0468	
Pima Indians	0.36309+-4.3776e- 002	0.69672+- 7.4829e- 002	
Soanr	0.61524+-7.1800e- 002	0.85024+-6.0008e- 002	
Heart-statlog	.57407+-8.4186e- 002	0.67778+-6.0607e- 002	

#### We can extend this method to Nonlinear Classifier

One of the hardest parts of writing a research paper can be just finding a good topic to write about.

Some ideas:

- 1. Finding a new method to separate data sets
- 2. New efficient optimization model for the previous ideas.
- 3. Solving the existence ideas with a new method.
- 4. Extending the currents methods for binary classification to

Multi-class classification

# Resources: Datasets

- UCI Repository: <u>http://www.ics.uci.edu/~mlearn/MLRepository.htm</u>
- UCI KDD Archive:

http://kdd.ics.uci.edu/summary.data.application. html

- Statlib: <u>http://lib.stat.cmu.edu/</u>
- Delve: <u>http://www.cs.utoronto.ca/~delve/</u>

# **Resources:** Journals

- Journal of Machine Learning Research Machine Learning
- IEEE Transactions on Neural Networks
- IEEE Transactions on Pattern Analysis and Machine Intelligence
- Annals of Statistics
- Journal of the American Statistical Association

• ...

# Resources: Conferences

- International Conference on Machine Learning (ICML)
- European Conference on Machine Learning (ECML)
- Neural Information Processing Systems (NIPS)
- Computational Learning
- International Joint Conference on Artificial Intelligence (IJCAI)
- ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)
- IEEE Int. Conf. on Data Mining (ICDM)

## Questions?




## Thanks for your **Atention**

## Charles Bridge Prague

