Large Scale Support Vector Machines

Huanle Xu

Motivation

- Introduce the widely used classification tool: Support Vector Machine (SVM)
- Understand the model and parameter estimation method in terms of big data

Motivation

Suppose we have 50 photographs of elephants and 50 photos of tigers.





We digitize them into 100 x 100 pixel images, so we have $x \in \mathbb{R}^n$ where n = 10,000.

VS.

Now, given a new (different) photograph we want to answer the question: **is it an elephant or a tiger?** [we assume it is one or the other.]



photos, how to make the SVM training scalable?

SVMs: History

- SVMs introduced in COLT-92 by Boser, Guyon
 & Vapnik. Became rather popular since.
- Theoretically well motivated algorithm: developed from Statistical Learning Theory (Vapnik & Chervonenkis) since the 60s.
- Empirically good performance: successful applications in many fields (bioinformatics, text, image recognition, . . .)

SVMs: History

- Centralized website: www.kernelmachines.org.
- Several textbooks, e.g. "An introduction to Support Vector Machines" by Cristianini and Shawe-Taylor is one.
- A large and diverse community work on them: from machine learning, optimization, statistics, neural networks, functional analysis, etc.

Linear SVMs

- Data
 - Training examples: $(x_1, y_1), \ldots, (x_n, y_n)$
 - Each $x_i \in \mathbb{R}^d, y_i \in \{+1, -1\}$
 - Want to find a hyperplane $y = w^T x + b$ to separate "+" from "-"
- What's the best hyperplane defined by w ?



- Distance from the separating hyperplance corresponds to the "confidence" of prediction
- Example: We have more confidence to say A and B belong to "+" than C



- Support Vectors: Examples closest to the hyperplane
- Margin *P* : width of separation between support vectors of classes.



Support vector

• Distance from example to the separator is :

$$r = y \frac{w^T x + b}{\|w\|}$$

• Proof:

 $\begin{aligned} x' - x//w, \text{ unit vector is } w/\|w\|, \\ \text{so line is } rw/\|w\|, \ x' = x - yrw/\|w\| \\ \text{since } x' \text{ is on the separator, } w^T x' + b = 0 \\ \text{so } w^T (x - yrw/\|w\|) = 0, \ w = \sqrt(w^T w), \\ \text{so } w^T x - yr\|w\| + b = 0, \\ \text{then we get } r = y \frac{w^T x + b}{\|w\|} \end{aligned}$



- Assume that all data is at least distance 1 from the hyperplane, then the following constraints follow for a training set $\{(x_i, y_i)\}_{i=1}^n$ $y_i(w^T x_i + b) > 1$
- For support vectors, the inequality becomes an equality
- Recall that $r = y \frac{w^T x + b}{\|w\|}$ Margin is: $\rho = \frac{2}{\|w\|}$

Linear SVMs

- Note that we assume that all data points are linearly separated by the hyperplane.
- The margin is invariant to scaling of parameters.
 - i.e. by changing w, b to 5w, 5b, the margin doesn't change

Linear SVMs

• Maximize the margin

Good according to intuition, theory (VC dimension) & practice

• The problem of linear SVMs is formulated as:

$$\max_{w} \rho = \frac{2}{\|w\|}$$

s.t. $y_i(w^T x_i + b) \ge 1 \quad \forall i = 1, \dots, n$

• An equivalent form is:

$$\min_{w} \frac{1}{2} \|w\|^2$$

s.t. $y_i(w^T x_i + b) \ge 1 \quad \forall i = 1, \dots, n$

Non-Linear Separable SVMs

- In reality, training samples are usually not linearly separable.
- Soft Margin Classification
 - Idea: allow errors but introduce slack variable ξ_i to penalize errors
 - Still try to minimize training set errors, and to place hyperplane "far" from each class (large margin)



Soft Margin Classification

- The problem becomes:
 - $\min_{w} \frac{1}{2} \|w\|^{2} + C \sum_{i} \xi_{i}$ s.t. $y_{i}(w^{T}x_{i} + b) \ge 1 \xi_{i}, \quad \xi_{i} \ge 0 \quad \forall i = 1, \dots, n$
 - Minimize $||w||^2$ plus the number of training mistakes
 - Set C using cross validation

Soft Margin Classification

- If point x_i is on the wrong side of the margin then get penalty ξ_i
- Thus all mistakes are not equally bad!



For each datapoint: If margin \ge 1, don't care If margin < 1, pay linear penalty

Slack Penalty C

$$\min_{w} \frac{1}{2} \|w\|^{2} + C \sum_{i} \xi_{i}$$

s.t. $y_{i}(w^{T}x_{i} + b) \ge 1 - \xi_{i}, \quad \xi_{i} \ge 0 \quad \forall i = 1, \dots, n$

- What is the role of penalty C:
 - $C \equiv 0$: can set ξ_i to anything, then w=0 (basically ignore the data)
 - $-C = \infty$: Only want w,b to separate the data (0,0)



Soft Margin Classification

- SVM in the "natural" form $\underset{w,b}{\operatorname{arg\,min}} \frac{1}{2} \|w\|^{2} + C \sum_{i=1}^{n} \max\{0, 1 - y_{i}(w^{T}x_{i} + b)\}$ Margin Regularization Parameter Empirical loss L
- SVM uses "Hinge Loss":

Non-linear Separable SVMs

- Linear classifiers aren't complex enough sometimes.
 - Map data into a richer feature space including non-linear features
 - Then construct a hyperplane in that space so all other equations are the same



Non-linear Separable SVMs

• Formally, process the data with:

 $x \mapsto \Phi(x)$

• Then learn the map from $\Phi(x)$ to y

$$f(x) = w \cdot \Phi(x) + b$$



Example: Polynomial Mapping

 $\Phi: R^2 \to R^3$ $(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{(2)} x_1 x_2, x_2^2)$



Example: MNIST

- Data: 60,000 training examples, 10000 test examples, 28x28
- Linear SVM has around 8.5% test error. Polynomial SVM has around 1% test error.



MINST Results

Classifier	Test Error
linear	8.4%
3-nearest-neighbor	2.4%
RBF-SVM	1.4 %
Tangent distance	1.1 %
LeNet	1.1 %
Boosted LeNet	0.7 %
Translation invariant SVM	0.56 %

Choosing a good mapping $\Phi(\cdot)$ (encoding prior knowledge + getting right complexity of function class) for your problem improves results.

SVM: How to Estimate w, b

• We take the soft margin classification for example:

$$\min_{w} \frac{1}{2} \|w\|^{2} + C \sum_{i} \xi_{i}$$
s.t. $y_{i}(w^{T}x_{i} + b) \ge 1 - \xi_{i}, \quad \xi_{i} \ge 0 \quad \forall i = 1, \dots, n$

• Standard way: Use a solver!

 Solver: software for finding solutions to "common" optimization problems, e.g. LIBSVM (<u>http://www.csie.ntu.edu.tw/~cjlin/libsvm/</u>)

• Problems: Solvers are inefficient for big data!

SVM: How to Estimate w, b

- Want to estimate w,b !
- Alternative approach: $s.t. \forall i \ y_i(w^T x_i + b) \ge 1 \xi_i, \ \xi_i \ge 0$ - Want to minimize **f(w,b)**

 $\min_{w} \frac{1}{2} \|w\|^2 + C \sum \xi_i$

$$f(w,b) = \frac{1}{2} \sum_{j=1}^{d} (w^{(j)})^2 + C \sum_{i=1}^{n} \max\{0, 1 - y_i (\sum_{j=1}^{d} w^{(j)} x_i^{(j)} + b)\}$$

- How to minimize convex functions f(z)
- Use gradient descent: $\min_{z} f(z)$
- Iterate: $z_{t+1} \leftarrow z_t \eta f'(z_t)$



SVM: How to Estimate w?

• Want to minimize *f(w,b)*:

$$f(w,b) = \frac{1}{2} \sum_{j=1}^{d} (w^{(j)})^2 + C \sum_{i=1}^{n} \max\{0, 1 - y_i(\sum_{j=1}^{d} w^{(j)} x_i^{(j)} + b)\}$$

Empirical loss L

• Compute the gradient $\nabla(j) w.r.t w^{(j)}$

$$\nabla(j) = \frac{\partial f(w,b)}{\partial w^{(j)}} = w^{(j)} + C \sum_{i=1}^{n} \frac{\partial L(x_i, y_j)}{\partial w^{(j)}}$$

 $\frac{\partial L(x_i, y_j)}{\partial w^{(j)}} = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \ge 1\\ -y_i x_i^{(j)} & \text{otherwise} \end{cases}$

SVM: How to Estimate w?

• Gradient descent:

Iterate untial convergence:

• For j = 1, ..., d

- Evaluate:
$$\nabla(j) = \frac{\partial f(w,b)}{\partial w^{(j)}} = w^j + C \sum_{i=1}^n \frac{\partial L(x_i,y_i)}{\partial w^{(j)}}$$

- Update: $w^{(j)} = w^{(j)} - \eta \nabla(j)$
 $\eta \dots$ learning rate parameter
 $C \dots$ regularization parameter

- Problem:
 - Computing $\nabla(j)$ takes O(n) time
 - n ... size of the training dataset

SVM: How to Estimate w?

- Stochastic Gradient Descent We just had: $\nabla(j) = w^{(j)} + C \sum_{i=1}^{n} \frac{\partial L(x_i, y_i)}{\partial w^{(j)}}$
 - Instead of evaluating gradient over all examples, evaluate it for each individual training example

$$\nabla(j,i) = w^{(j)} + C \frac{\partial L(x_i, y_i)}{\partial w^{(j)}}$$

• Stochastic gradient descent:

Iterate untial convergence:

• For
$$i = 1, \ldots, n$$

- For
$$j = 1, ..., d$$

* Evaluate: $\nabla(j, i)$
* Upadate: $w^{(j)} \leftarrow w^{(j)} - \eta \nabla(j, i)$



For optimizing *f(w,b)* within reasonable quality SGD-SVM is super fast

SGD vs. Batch Conjugate Gradient

• SGD on full dataset vs. Batch Conjugate

Gradient on a sample of *n* training examples



many times is better than doing a complicated (but slow) BCG update a few times



k... condition number

• Need to choose learning rate η and t_0

$$w_{t+1} \leftarrow w_t - \frac{\eta_t}{t+t_0} \left(w_t + C \frac{\partial L(x_i, y_i)}{\partial w} \right)$$

- Leon suggests:
 - Choose t_0 so that the expected initial updates are comparable with the expected size of the weights
 - Choose η :
 - Select a small subsample
 - Try various rates η (e.g., 10,1,0.1,0.01,...)
 - Pick the one that most reduces the cost
 - Use η for next 100k iterations on the full dataset

• Sparse Linear SVM:

- Feature vector x_i is sparse (contains many zeros)

- Do not do: $x_i = [0, 0, 0, 1, 0, 0, 0, 0, 5, 0, 0, 0, 0, 0, 0, ...]$
- But represent x_i as a sparse vector $\mathbf{x_i} = [(\mathbf{4}, \mathbf{1}), (\mathbf{9}, \mathbf{5}), \ldots]$
- Can we do the SGD update more efficiently?

$$w \leftarrow w - \eta \left(w + C \frac{\partial L(x_i, y_i)}{\partial w} \right)$$

– Approximated in 2 steps:

$$w \leftarrow w - \eta C \frac{\partial L(x_i, y_i)}{\partial w}$$
$$w \leftarrow w(1 - \eta)$$

Cheap: x_i is sparse and so few coordinates **j** of **w** will be updates Expensive: **w** is not sparse, all coordinates need to be updated

- Solution 1: $\mathbf{w} = \mathbf{s} \cdot \mathbf{v}$
 - Represent vector w as the product of scalar s and the vector v
 - Then the update procedure is:
 - 1) $v = v \eta C \frac{\partial L(x_i, y_i)}{\partial w}$ Two step update procedure: • 2) $s = s(1 - \eta)$ 1. $w \leftarrow w - \eta C \frac{\partial L(x_i, y_i)}{\partial w}$
- Solution 2:

2. $w \leftarrow w(1 - \eta)$

- Perform only step 1) for each training example
- Perform step 2) with lower frequency and higher η

• Stopping criteria:

How many iterations of SGD?

- Early stopping with cross validation
 - Create validation set
 - Monitor cost function on the validation set
 - Stop when loss stops decreasing

• Stopping criteria:

How many iterations of SGD?

- Early Stopping
 - Extract two disjoint subsamples **A** and **B** of training data
 - Train on **A**, stop by validating on **B**
 - Number of epochs is an estimate of **k**
 - Train for **k** epochs on the full dataset

What about Multiple Classes?

- Idea 1:

 One against all
 Learn 3 classifiers
 - + vs. {o,-}
 - - vs. {o,+}
 - o vs. {+,-}

Obtain: $w_{+}b_{+}, w_{-}b_{-}, w_{o}b_{o}$

Return class c

 $\operatorname{arg\,max}_{c} w_{c} x + b_{c}$



What about Multiple Classes?

- Idea 2:
 - Learn 3 sets of weights simultaneously
 - Want the correct class to have highest margin:

$$w_{y_i}x_i + b_{y_i} \ge 1 + w_c x_i + b_c \ \forall c \neq y_i, \forall i$$



Multiclass SVM

- Optimization problem:
 - $\min_{w,b} \frac{1}{2} \sum_{c} ||w_{c}||^{2} + C \sum_{i=1}^{n} \xi_{i}$ $w_{y_{i}} x_{i} + b_{y_{i}} \ge w_{c} x_{i} + b_{c} + 1 - \xi_{i} \,\forall c \neq y_{i}, \xi_{i} \ge 0, \forall i$
 - To obtain parameters w_c, b_c for each class c, we can use similar techniques as for 2 class SVM
- SVM is widely perceived a very powerful learning algorithm

Reference

- http://www.stanford.edu/class/cs246/slides/13-svm.pdf
- <u>http://www.stanford.edu/class/cs276/handouts/lecture14-</u>
 <u>SVMs.ppt</u>
- http://i.stanford.edu/~ullman/pub/ch12.pdf
- http://www.svms.org/tutorials/
- <u>http://www.cs.columbia.edu/~kathy/cs4701/documents/jaso</u>
 <u>n_svm_tutorial.pdf</u>
- http://www.csie.ntu.edu.tw/~cjlin/libsvm/
- Chang, E, Zhu, K, Wang, H, Bai, H, Li, J, Qiu, Z, and Cui, H. PSVM: Parallelizing support vector machines on distributed computers. NIPS, 20:257-264. 2007.

In-class Practice



 Consider building an SVM over the (very little) data set shown in above figure, compute the e SVM decision boundary.