#### Transduction

## Outline

- Introduction
- History
- Theoretical Details
- Algorithm Walkthrough
- Implementation Results
- Strengths and Weaknesses
- Bells and Whistles
- Current Trajectory
- Conclusion

# **Introduction --- The Problem Setting**

Given:

- Labeled set
- Unlabeled set
- Labels \in {-1, 1}

Find set of labels that best fits unlabeled set



# **Introduction --- The Problem Setting**

Given:

- Labeled set
- Unlabeled set
- Labels \in {-1, 1}

Find set of labels that best fits unlabeled set

Note that we aren't extending to unseen examples!



# **Introduction --- Why we might care**

Given:

- Labeled set
- Unlabeled set
- Labels \in {-1, 1}

Find set of labels that best fits unlabeled set

Note that we aren't extending to unseen examples!

Theoretically:

- For problems with partially labeled data
- When we don't need a general function

# **Introduction --- Why we might care**

Concretely:

- Text classification (Joachims, 1999)
- Sentiment analysis (Talukdar and Crammer, 2009)
- Image classification (Yu et al., 2012)

Theoretically:

- For problems with partially labeled data
- When we don't need a general function

## History

1963 --- Maximum-margin hyperplane algorithm by Vapnik

1971 --- Vapnik-Chervonenkis dimension

1992 --- Kernel trick by Boser, Guyon, and Vapnik

1998 --- Transduction by Vapnik



http://www.lecun.org/gallery/libpro/20011121-allyourbayes/dsc01228-02-h.jpg

Main Points:

- We want our model to have the least risk possible
- Getting the tightest bound on the risk is best, theoretically
- Transduction theoretically follows the tightest bound

Main Points:

- We want our model to have the least risk possible
- Getting the tightest bound on the risk is best, theoretically
- Transduction theoretically follows the tightest bound

Admissions of Ignorance:

- I'm not exactly sure what the function f is
- I don't know why transduction theoretically follows the tightest bound

Useful tidbit:

• VC = Vapnik-Chervonenkis

True Risk Functional (accounting for the true probability distribution):

$$R(\alpha) = \int |y - f(x, \alpha)| dP(x, y)$$

Why can't we calculate this very easily?

 We would need to know P(x,y) beforehand; but if we knew P(x,y), then we wouldn't need to do machine learning (marginalize y to get P(x), divide P(x,y) by P(x) to get P(y|x), choose most probably y)

True Risk Functional (accounting for the true probability distribution):

$$R(\alpha) = \int |y - f(x, \alpha)| dP(x, y)$$

Empirical Risk Functional (making do with the data we have):

$$R_{emp}(\alpha) = \sum_{i=1}^{\ell} |y_i - f(x_i, \alpha)|$$



VC Entropy Bound

$$R(\alpha) \le R_{emp}(\alpha) + \sqrt{\frac{H_P^{\Lambda}(2\ell) - \ln \eta}{\ell}}$$

Growth Function Bound

$$R(\alpha) \le R_{emp}(\alpha) + \sqrt{\frac{\ln G^{\Lambda}(2\ell) - \ln \eta}{\ell}}$$

VC Dimension Bound

$$R(\alpha) \le R_{emp}(\alpha) + \sqrt{\frac{h(\ln\frac{2\ell}{h}+1) - \ln\eta}{\ell}}$$

Again, I'm not sure how these theoretical details relate to the implementation of transduction.

I have put them in the slides because Vapnik seems to think that they are important.

He does talk about structural risk minimization and applies the bounds to come up with an inequality over test set vectors, which seems like something to do with transduction, but nearly all of that went over my head.

OP 2 (Transductive SVM (non-sep. case)) Minimize over  $(y_1^*, ..., y_n^*, \vec{w}, b, \xi_1, ..., \xi_n, \xi_1^*, ..., \xi_k^*)$ :  $\frac{1}{2}||\vec{w}||^2 + C\sum_{i=0}^n \xi_i + C^* \sum_{j=0}^k \xi_j^*$  $\forall_{i=1}^{n}: y_{i}[\vec{w} \cdot \vec{x}_{i} + b] > 1 - \xi_{i}$ subject to:  $\forall_{i=1}^{k} : y_{i}^{*}[\vec{w} \cdot \vec{x}_{i}^{*} + b] \ge 1 - \xi_{i}^{*}$  $\forall_{i-1}^n : \xi_i > 0$  $\forall_{i=1}^k : \xi_i^* > 0$ 

Algorithm TSVM:

Input: - training examples  $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ - test examples  $\vec{x}_1^*, \dots, \vec{x}_k^*$  $-C,C^*$ : parameters from OP(2) Parameters:  $-num_+$ : number of test examples to be assigned to class + - predicted labels of the test examples  $y_1^*, \ldots, y_k^*$ Output:  $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ Classify the test examples using  $\langle \vec{w}, b \rangle$ . The  $num_+$  test examples with the highest value of  $\vec{w} * \vec{x}_i^* + b$  are assigned to the class +  $(y_i^* := 1)$ ; the remaining test examples are assigned to class -  $(y_i^* := -1)$ .  $C_{-}^{*} := 10^{-5}$ : // some small number  $C_{+}^{*} := 10^{-5} * \frac{num_{+}}{k - num_{+}};$ while  $((C_{-}^* < C^*) \parallel (C_{+}^* < C^*))$  { // Loop 1  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*)...(\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ while  $(\exists m, l: (y_m^* * y_l^* < 0) \& (\xi_m^* > 0) \& (\xi_l^* > 0) \& (\xi_m^* + \xi_l^* > 2))$  { // Loop 2 // take a positive and a negative test
// example, switch their labels, and retrain  $y_m^* := -y_m^*;$  $y_l^* := -y_l^*;$  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1) ... (\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*) ... (\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$  $C_{+}^{*} := \min(C_{-}^{*} * 2, C^{*});$   $C_{+}^{*} := \min(C_{+}^{*} * 2, C^{*});$ return $(y_1^*, ..., y_k^*)$ ;

Input:

- training examples  $(\vec{x}_1, y_1), ..., (\vec{x}_n, y_n)$
- test examples  $\vec{x}_1^*, ..., \vec{x}_k^*$

Parameters:

- $-C, C^*$ : parameters from OP(2)
  - $-num_+$ : number of test examples to be assigned to class +

Output:

- predicted labels of the test examples  $y_1^*, \ldots, y_k^*$ 

Note that Input and Output follows the definition of a transduction problem

Input:

- training examples  $(\vec{x}_1, y_1), ..., (\vec{x}_n, y_n)$
- test examples  $\vec{x}_1^*, ..., \vec{x}_k^*$

Parameters:

- $-C,C^*$ : parameters from OP(2)
  - $-num_+$ : number of test examples to be assigned to class +

Output:

– predicted labels of the test examples  $y_1^*, ..., y_k^*$ 

As for Parameters:

OP 2 (Transductive SVM (non-sep. case)) Minimize over  $(y_1^*, ..., y_n^*, \vec{w}, b, \xi_1, ..., \xi_n, \xi_1^*, ..., \xi_k^*)$ :

$$\begin{aligned} \frac{1}{2} ||\vec{w}||^2 + C \sum_{i=0}^n \xi_i + C^* \sum_{j=0}^k \xi_j^* \\ subject \ to: \qquad \forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x}_i + b] \ge 1 - \xi_i \\ \forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x}_j^* + b] \ge 1 - \xi_j^* \\ \forall_{i=1}^n : \xi_i > 0 \\ \forall_{j=1}^k : \xi_j^* > 0 \end{aligned}$$

 $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ 

The function  $solve\_svm\_qp$  refers to quadratic programs of the following type.

```
OP 3 (Inductive SVM (primal))

Minimize over (\vec{w}, b, \vec{\xi}, \vec{\xi^*}):

\frac{1}{2} ||\vec{w}||^2 + C \sum_{i=1}^n \xi_i + C_-^* \sum_{j:y_j^*=-1} \xi_j^* + C_+^* \sum_{j:y_j^*=1} \xi_j^*
subject to: \forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x_i} + b] \ge 1 - \xi_i

\forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x_j} + b] \ge 1 - \xi_j^*
```

 $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ 

The function  $solve\_svm\_qp$  refers to quadratic programs of the following type.

**OP 3 (Inductive SVM (primal))** Minimize over  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*})$ :

$$\frac{1}{2} ||\vec{w}||^2 + C \sum_{i=1}^n \xi_i + C_-^* \sum_{j:y_j^* = -1} \xi_j^* + C_+^* \sum_{j:y_j^* = 1} \xi_j^*$$
  
subject to:  $\forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x_i} + b] \ge 1 - \xi_i$   
 $\forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x_j} + b] \ge 1 - \xi_j^*$ 

Why aren't these constraints part of OP 3?  $\forall_{i=1}^{n}: \xi_i > 0$  $\forall_{j=1}^{k}: \xi_j^* > 0$ 

 $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ 

The function  $solve\_svm\_qp$  refers to quadratic programs of the following type.

**OP 3 (Inductive SVM (primal))** Minimize over  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*})$ :

$$\frac{1}{2} ||\vec{w}||^2 + C \sum_{i=1}^n \xi_i + C_-^* \sum_{j:y_j^* = -1} \xi_j^* + C_+^* \sum_{j:y_j^* = 1} \xi_j^*$$
  
subject to:  $\forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x_i} + b] \ge 1 - \xi_i$   
 $\forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x_j} + b] \ge 1 - \xi_j^*$ 

Why aren't these constraints part of OP 3?  $\forall_{i=1}^{n} : \xi_i > 0$  $\forall_{j=1}^{k} : \xi_j^* > 0$ 

Ideally, they are fulfilled by the other constraints

 $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ 

The function  $solve\_svm\_qp$  refers to quadratic programs of the following type.

**OP 3 (Inductive SVM (primal))** Minimize over  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*})$ :

$$\frac{1}{2} ||\vec{w}||^2 + C \sum_{i=1}^n \xi_i + C_{-}^* \sum_{j:y_j^* = -1} \xi_j^* + C_{+}^* \sum_{j:y_j^* = 1} \xi_j^*$$
  
subject to:  $\forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x_i} + b] \ge 1 - \xi_i$   
 $\forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x_j} + b] \ge 1 - \xi_j^*$ 

Why aren't these constraints part of OP 3?  $\forall_{i=1}^{n} : \xi_i > 0$  $\forall_{j=1}^{k} : \xi_j^* > 0$ 

Ideally, they are fulfilled by the other constraints

Then why are they part of OP 2? I don't know...

Classify the test examples using  $\langle \vec{w}, b \rangle$ . The  $num_+$  test examples with the highest value of  $\vec{w} * \vec{x}_j^* + b$  are assigned to the class  $+ (y_j^* := 1)$ ; the remaining test examples are assigned to class  $- (y_j^* := -1)$ .  $C_-^* := 10^{-5}$ ; // some small number  $C_+^* := 10^{-5} * \frac{num_+}{k-num_+}$ ;

Classification is based on a model trained on only labeled data

Classify the test examples using  $\langle \vec{w}, b \rangle$ . The  $num_+$  test examples with the highest value of  $\vec{w} * \vec{x}_j^* + b$  are assigned to the class  $+ (y_j^* := 1)$ ; the remaining test examples are assigned to class  $- (y_j^* := -1)$ .  $C_-^* := 10^{-5};$  // some small number  $C_+^* := 10^{-5} * \frac{num_+}{k-num_+};$ 

The small numbers refer back to OP 3

The function  $solve\_svm\_qp$  refers to quadratic programs of the following type.

**OP 3 (Inductive SVM (primal))** Minimize over  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*})$ :

 $\frac{1}{2} ||\vec{w}||^2 + C \sum_{i=1}^n \xi_i + C_-^* \sum_{j:y_j^* = -1} \xi_j^* + C_+^* \sum_{j:y_j^* = 1} \xi_j^*$ subject to:  $\forall_{i=1}^n : y_i [\vec{w} \cdot \vec{x_i} + b] \ge 1 - \xi_i$  $\forall_{j=1}^k : y_j^* [\vec{w} \cdot \vec{x_j} + b] \ge 1 - \xi_j^*$ 

while  $((C_{-}^* < C^*) \parallel (C_{+}^* < C^*))$ // Loop 1  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1) ... (\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*) ... (\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ while  $(\exists m, l: (y_m^* * y_l^* < 0) \& (\xi_m^* > 0) \& (\xi_l^* > 0) \& (\xi_m^* + \xi_l^* > 2))$  { // Loop 2 // take a positive and a negative test  $y_{m}^{*} := -y_{m}^{*};$  $y_1^* := -y_1^*$ ; // example, switch their labels, and retrain  $(\vec{w}, b, \vec{\xi}, \vec{\xi}^*) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*)...(\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$  $C_{-}^{*} := \min(C_{-}^{*} * 2, C^{*});$   $C_{+}^{*} := \min(C_{+}^{*} * 2, C^{*});$ 

Loop 1: While the small numbers are small, keep changing them

while  $((C_{-}^* < C^*) \parallel (C_{+}^* < C^*))$  { // Loop 1  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1) ... (\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*) ... (\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ while  $(\exists m, l: (y_m^* * y_l^* < 0) \& (\xi_m^* > 0) \& (\xi_l^* > 0) \& (\xi_m^* + \xi_l^* > 2))$  { // Loop 2 // take a positive and a negative test  $y_m^* := -y_m^*;$  $y_{i}^{*} := -y_{i}^{*}$ : // example, switch their labels, and retrain  $(\vec{w}, b, \vec{\xi}, \vec{\xi}^*) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*)...(\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ }  $C_{-}^{*} := min(C_{-}^{*} * 2, C^{*});$  $C_{+}^{*} := min(C_{+}^{*} * 2, C^{*});$ 

# Loop 2: Swap test instance labels while corresponding penalties don't match

Algorithm TSVM:

Input: - training examples  $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ - test examples  $\vec{x}_1^*, \dots, \vec{x}_k^*$  $-C,C^*$ : parameters from OP(2) Parameters:  $-num_+$ : number of test examples to be assigned to class + - predicted labels of the test examples  $y_1^*, \ldots, y_k^*$ Output:  $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ Classify the test examples using  $\langle \vec{w}, b \rangle$ . The  $num_+$  test examples with the highest value of  $\vec{w} * \vec{x}_i^* + b$  are assigned to the class +  $(y_i^* := 1)$ ; the remaining test examples are assigned to class -  $(y_i^* := -1)$ .  $C_{-}^{*} := 10^{-5}$ : // some small number  $C_{+}^{*} := 10^{-5} * \frac{num_{+}}{k - num_{+}};$ while  $((C_{-}^* < C^*) \parallel (C_{+}^* < C^*))$  { // Loop 1  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*)...(\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ while  $(\exists m, l: (y_m^* * y_l^* < 0) \& (\xi_m^* > 0) \& (\xi_l^* > 0) \& (\xi_m^* + \xi_l^* > 2))$  { // Loop 2 // take a positive and a negative test
// example, switch their labels, and retrain  $y_m^* := -y_m^*;$  $y_l^* := -y_l^*;$  $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1) ... (\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*) ... (\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$  $C_{+}^{*} := \min(C_{-}^{*} * 2, C^{*});$   $C_{+}^{*} := \min(C_{+}^{*} * 2, C^{*});$ return $(y_1^*, ..., y_k^*)$ ;

Red is 1, Blue is -1

+ is given 1, - is given -1

Circles are test data

Dashed line is initial hyperplane

Solid line is final hyperplane



0	0.55	1
0.5	0.49	1
0	0.33	?
0.3	0.4	?
0.6	0.9	?
0.8	0.6	?
0.4	0.51	-1
0.45	0.3	-1



```
Initial
Initial
predictions: [1,-1,1,-1]
weights: [-0.3499999989415537,0.2299999993580931]
bias: 0.011620685479047073
training penalties:
[0.8618793073756819, 1.0506793068090083, 0.9889206780774836, 0.9231206782648367]
testing penalties: Any[]
#----
C star plus: 9.999999999999999-6
```

 $(\vec{w}, b, \vec{\xi}, \_) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [], C, 0, 0);$ 

```
Initial
Initial
predictions: [1,-1,1,-1]
weights: [-0.3499999989415537,0.2299999993580931]
bias: 0.011620685479047073
training penalties:
[0.8618793073756819, 1.0506793068090083, 0.9889206780774836, 0.9231206782648367]
testing penalties: Any[]
#-----
C star plus: 9.999999999999999-6
  C_star_minus: 9.999999999999999-6
```

Classify the test examples using  $\langle \vec{w}, b \rangle$ . The  $num_+$  test examples with the highest value of  $\vec{w} * \vec{x}_j^* + b$  are assigned to the class  $+ (y_j^* := 1)$ ; the remaining test examples are assigned to class  $- (y_j^* := -1)$ .

```
Initial
Initial
predictions: [1,-1,1,-1]
weights: [-0.3499999989415537,0.2299999993580931]
bias: 0.011620685479047073
training penalties:
[0.8618793073756819, 1.0506793068090083, 0.9889206780774836, 0.9231206782648367]
testing penalties: Any[]
#-----
C star plus: 9.999999999999999-6
  C_star_minus: 9.999999999999999-6
```

// some small number

 $C_{-}^{*} := 10^{-5};$  $C_{+}^{*} := 10^{-5} * \frac{num_{+}}{k - num_{+}};$ 

```
Iteration: 1
Iteration: 1
predictions: [1,-1,1,-1]
weights: [-0.35000499744602587,0.23000229886942783]
bias: 0.04813871789811412
training penalties:
[0.8253600102296041,1.0141626468847826,1.0254378838490155,0.9596371512141343]
testing penalties:
[0.8762113078710513, 1.0353887820662886, 0.955112927968292, 0.9063867484489777]
           ----#
#----
C_star_plus: 1.9999999999999998e-5
  C star minus: 1.9999999999999998e-5
```

 $(\vec{w}, b, \vec{\xi}, \vec{\xi^*}) := solve\_svm\_qp([(\vec{x}_1, y_1)...(\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*)...(\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*);$ 

```
Iteration: 1
Iteration: 1
predictions: [1,-1,1,-1]
weights: [-0.35000499744602587,0.23000229886942783]
bias: 0.04813871789811412
training penalties:
[0.8253600102296041,1.0141626468847826,1.0254378838490155,0.9596371512141343]
testing penalties:
[0.8762113078710513, 1.0353887820662886, 0.955112927968292, 0.9063867484489777]
           -----#
#----
C star plus: 1.9999999999999998e-5
  C star minus: 1.9999999999999998e-5
```

 $C_{-}^{*} := \min(C_{-}^{*} * 2, C^{*});$  $C_{+}^{*} := \min(C_{+}^{*} * 2, C^{*});$ 

$$\texttt{while}\left(\exists m,l: (y_m^* * y_l^* < 0) \& (\xi_m^* > 0) \& (\xi_l^* > 0) \& (\xi_m^* + \xi_l^* > 2) \right) \ \{ // \ \text{Loop } 2 \}$$

 $\begin{array}{ll} y_m^* := -y_m^*; & // \text{ take a positive and a negative test} \\ y_l^* := -y_l^*; & // \text{ example, switch their labels, and retrain} \\ (\vec{w}, b, \vec{\xi}, \vec{\xi}^*) := solve\_svm\_qp([(\vec{x}_1, y_1) ... (\vec{x}_n, y_n)], [(\vec{x}_1^*, y_1^*) ... (\vec{x}_k^*, y_k^*)], C, C_-^*, C_+^*); \end{array}$
```
Iteration: 16
Iteration: 16
predictions: [1,-1,1,-1]
weights: [-0.5138399971676509,0.3053663988809765]
bias: 0.14047075283260615
training penalties:
[0.6915777202844386,0.9668197028022565,1.0906716098973297,1.0008526662735882]
testing penalties:
[0.7587583331474819,1.1084653108507676,0.8930034840908901,0.9126185920347551]
#### Sum: 2.001468794941658
****
   index1: 2 1.1084653108507676
   index2: 3 0.8930034840908901
I*****************************
```

```
Iteration: 17
Swapped 2 and 3
Iteration: 17
Swapped 2 and 3
predictions: [1,1,-1,-1]
weights: [-0.7104479950209311,-0.02231359821962453]
bias: 0.25134504027522153
training penalties:
[0.7609274312467963,1.1148126128656513,0.9557858996772918,0.9249493555522572]
testing penalties:
[0.7560184447407811,0.9707147951327828,0.8049940024740828,0.6695984829281191]
           -----#
C_star_plus: 0.655359999999999
   C star minus: 0.6553599999999999
```

```
Iteration: 18
Iteration: 18
predictions: [1,1,-1,-1]
weights: [-1.070895993602519,-0.274627197579686]
bias: 0.5010333531255504
training penalties:
[0.6500115980434442,1.1689819629926717,0.9326150774213058,0.9367419892329326]
testing penalties:
[0.5895936158842809,0.9300863178018592,0.611331272953632,0.4795402335037401]
#----
C_star_plus: 1.0
  C_star_minus: 1.0
```

```
while((C_{-}^* < C^*) \parallel (C_{+}^* < C^*)){
```

// Loop 1

0	0.55	1
0.5	0.49	1
0	0.33	1
0.3	0.4	1
0.6	0.9	-1
0.8	0.6	-1
0.4	0.51	-1
0.45	0.3	-1









































0	0.55	1
0.5	0.49	1
0	0.33	1
0.3	0.4	1
0.6	0.9	-1
0.8	0.6	-1
0.4	0.51	-1
0.45	0.3	-1



Experimental setup:

- C = 1.0 = penalty for training example margins
- C^{\*} = 1.0 = penalty for test example margins
- num\_{+} = percentage of positive examples in training set multiplied by number of examples in test set = expected number of positive examples in test set

Experimental setup (continued):

- 10 and 4 refer to 10-fold and 4-fold cross validation
- small refers to random 10% as training set, averaged over 10 runs
- linear svm and rbf svm are from sklearn (Python)

Successful Dataset:

- forests = forest type mapping
  - <u>https://archive.ics.uci.edu/ml/datasets/Forest+type+mapping</u>
  - 27 attributes
  - 523 instances

Swapping step got stuck going back and forth (bad solver?):

- vert = vertebral column
  - o <u>https://archive.ics.uci.edu/ml/datasets/Vertebral+Column</u>
  - 6 attributes
  - 310 instances
- hiv = HIV-1 protease cleavage
  - <u>https://archive.ics.uci.edu/ml/datasets/HIV-1+protease+cleavage</u>
  - 1 attribute (8 amino acids)
  - 6590 instances

	101 2515								
	forests 10	forests 4	small	vert 10	vert 4	vert small			
tsvm	0.90801	0.91202	0.84968	0.83871	0.84514	0.81362			
julia svm	0.89304	0.90218	0.85499	0.83871	0.83224	0.81433			
linear svm	0.92184	0.91214	0.86879	0.84194	0.75262	0.78996			
rbf svm	0.69599	0.69599	0.69639	0.67742	0.67741	0.6362			

foracta

What I learned:

What I learned: Don't trust theories

**Theorem 2** Algorithm 1 converges in a finite number of steps.

**Proof:** To prove this, it is necessary to show that loop 2 is exited after a finite number of iterations. This

## **Strengths and Weaknesses of TSVM**

Strengths:

• Explicitly accounts for test set examples

Weaknesses:

- Takes a long time to train
- Requires really good solver

## **Bells and Whistles**

- Train dual instead of primal (SVMLight)
- Adjust parameters (C, C^{\*}, num\_{+}) to find best settings
- Kernel Trick

### **Current Trajectory**

- Latest paper I could find was from 2013:
  - Bresson et al., <u>Multi-class Transductive Learning Based on { 1 Relaxations</u> of Cheeger Cut and Mumford-Shah-Potts Model, Journal of Mathematical Imaging and Vision (2013)
- Graph-based (reminiscent of Locally Linear Embedding)
  - Yu et al., <u>Adaptive Hypergraph Learning and its Application in Image</u> <u>Classification</u>, IEEE Transactions on Image Processing, vol 21: issue 7 (2012)
  - Zhang et al., <u>Transductive Learning on Adaptive Graphs</u>, AAAI 2010

### Conclusion

- Vapnik established fundamental theories for transduction
- Transductive Support Vector Machine is the first conscientious effort at transduction
- Current efforts in transduction seem to be in graph-based approaches
- Line between semi-supervised learning and transduction is very blurry

# Bibliography

- Support Vector Machine. Wikipedia. Accessed 17 Mar 2016.
- VC Dimension. Wikipedia. Accessed 17 Mar 2016.
- V. Vapnik. "Transductive Inference and Semi-Supervised Learning." In *Semi-supervised Learning*, 1st ed. Cambridge, Mass.: MIT Press, 2006, ch. 24.
- T. Joachims. "Transductive inference for text classification using support vector machines." ICML, 1999.
- P. Talukdar and K. Crammer. "New regularized algorithms for transductive learning." ECML PKDD, 2009.
- Yu et al. "Adaptive hypergraph learning and its application in image classification." IEEE Transactions on Image Processing, vol 21: issue 7, 2012.
- Bresson et al., "Multi-class Transductive Learning Based on & 1 Relaxations of Cheeger Cut and Mumford-Shah-Potts Model," Journal of Mathematical Imaging and Vision, 2013
- Zhang et al., "Transductive Learning on Adaptive Graphs," AAAI 2010