

Optimum reject options for multiclass classification

Johannes Kummert

October 19, 2015

Bachelor Thesis

Optimum reject options for multiclass classification

Johannes Kummert

1. Reviewer **Prof. Dr. Barbara Hammer**
Theoretical Computer Science
Bielefeld University

2. Reviewer **Benjamin Paaß**
Theoretical Computer Science
Bielefeld University

Supervisors Prof. Dr. Barbara Hammer and Benjamin Paaß

October 19, 2015

Johannes Kummert

Optimum reject options for multiclass classification

Bachelor Thesis, October 19, 2015

Reviewers: Prof. Dr. Barbara Hammer and Benjamin Paasßen

Supervisors: Prof. Dr. Barbara Hammer and Benjamin Paasßen

Abstract

Contents

1	Introduction	1
2	Reject Options	3
2.1	Two Classes	3
2.1.1	Reject Strategy	3
2.1.2	Optimal θ	3
2.1.3	Finding θ	4
2.2	Multi class Classification	5
2.2.1	Global Reject	7
2.2.2	Local Reject	7
2.2.3	Optimal Local Reject	8
2.2.4	Computation by Brute Force	8
2.2.5	Computation by Dynamic Programming	9
2.2.6	Greedy Computation	10
2.2.7	Evaluation	12
3	Application for SVM	13
4	Conclusions	15

Introduction

1

Reject Options

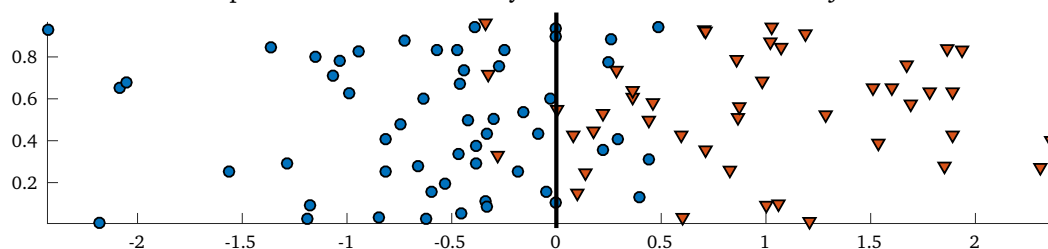
2.1 Two Classes

To make our way towards optimal rejects for multi class classification, we start of small by looking at a general two class classifier f that divides the space via a decision boundary.

$$f : \mathbb{R}^n \rightarrow \{1, 2\}$$

Let r be a measure of confidence that a point is part of its respective class, e.g. distance to the decision boundary. If $r(\bar{x})$ is large it means that \bar{x} is likely in the class it was assigned to.

Fig. 2.1: Example of two classes separated by a decision boundary. You can see falsely classified points near the boundary which we would like to reject.



2.1.1 Reject Strategy

We look for a threshold θ that defines for a given point \bar{x} whether it is rejected or not:

$$r(\bar{x}) < \theta : \bar{x} \text{ rejected}$$

2.1.2 Optimal θ

In order to find an optimal threshold we need to decide how to evaluate it. Optimally only wrongly classified data points are rejected (now referred to as true rejects). But generally there will be rejected although correctly classified points (false rejects). A

given labeled (by y) data set X is divided into the two sets L and E by applying f as a classifier with:

$$L = \{\bar{x} \in X \mid f(\bar{x}) = y(\bar{x})\} : \text{correctly classified points} \quad (2.1)$$

$$E = \{\bar{x} \in X \mid f(\bar{x}) \neq y(\bar{x})\} : \text{incorrectly classified points} \quad (2.2)$$

And by applying our reject strategy with the threshold θ , X is divided into A and R with:

$$A_\theta = \{\bar{x} \in X \mid r(\bar{x}) \geq \theta\} : \text{accepted points} \quad (2.3)$$

$$R_\theta = \{\bar{x} \in X \mid r(\bar{x}) < \theta\} : \text{rejected points} \quad (2.4)$$

The set T contains all true rejects and F all false rejects.

$$T_\theta = R \cap E \quad (2.5)$$

$$F_\theta = R \cap L \quad (2.6)$$

Naturally we want $|T_\theta|$ to be large and $|F_\theta|$ to be small. Since these two goals often contradict each other, e.g. more true rejects most times bring more false ones (see figure 2.2), there is no single optimal choice of θ in general. Still there is a set of values that we consider optimal. As shown above each θ corresponds to a tuple $(|T_\theta|, |F_\theta|)$. θ is optimal if

$$\nexists \theta' : |F_{\theta'}| \leq |F_\theta|, |T_{\theta'}| \geq |T_\theta|$$

and at least one term is unequal (TODO: wording?).

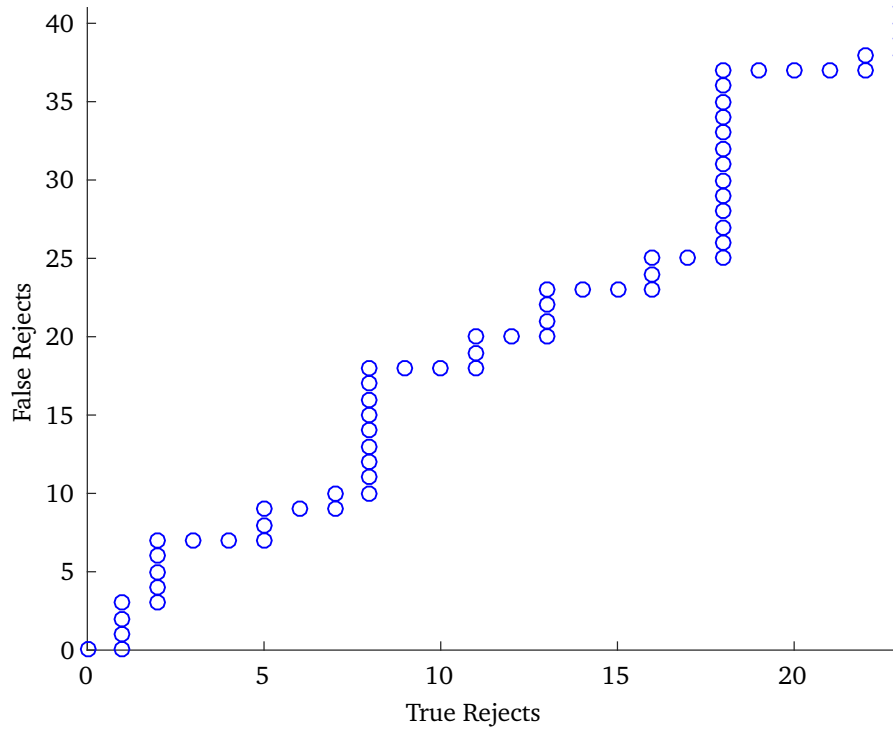
2.1.3 Finding θ

Now that we know how to evaluate our thresholds we still need to know where to look for them. Obviously θ needs to be in the range of r so

$$\Theta := \{\theta \in \mathbb{R} \mid \theta < \max_{\bar{x}} r(\bar{x})\}$$

is a set of possible thresholds. But it would leave us with an infinite search area (TODO: wording?) but this is easily reduced when recognized that it is unnecessary to consider multiple thresholds in between two points that are next to each other (when sorted according to r). So with $\Theta := \{r(\bar{x}) \mid \bar{x} \in X\}$ we have a feasible set of possible thresholds. It can be refined by looking again at our criteria for an optimal θ . We want $|T_\theta|$, the amount fo true rejects, to be large so we should skip cluster of true rejects since choosing a point in the middle of such cluster can not be better than the point at the end of it. Additionally we want $|F_\theta|$, the amount fo false rejects, to be small, so we should skip cluster of false rejects since a point in the middle of

Fig. 2.2: In this example we can see the corresponding true and false rejects if $r(\bar{x})$ is chosen as a threshold for each $\bar{x} \in X$. $r(\bar{x})$ is increasing from the bottom left corner to the upper right. As we can see, the number of false and true rejects are monotonically increasing with a more strict threshold.



such cluster can not be better than a point at the beginning of it. In conclusion this means that it is sufficient to consider only points at the beginning of clusters of false rejects, so

$$\Theta := \{r(\bar{x}_i) \mid \bar{x}_i \in L, \bar{x}_{i-1} \in E, \bar{x}_{i+1} \in L\}$$

where i is the index of the points in X when sorted according to r . This gives us an easily computed set of thresholds to consider to be optimal (see figures 2.4 and 2.3.)

2.2 Multi class Classification

Let us now expand the problem to a classifier for N multiple classes. By using a one vs all strategy, we get N binary classifiers f_i like in chapter 2.1 and an according measure of confidence r_i . Points are now classified in the class where confidence is maximal among all classes. This gives us a multi class classifier f with

$$f : \mathbb{R}^n \rightarrow \{1, \dots, i\} \quad (2.7)$$

$$\bar{x} \mapsto i \mid \arg \max_i r_i(\bar{x}) \quad (2.8)$$

Fig. 2.3: This figure shows the thresholds from before (see figure 2.2). The green marked points correspond to threshold we consider optimal.

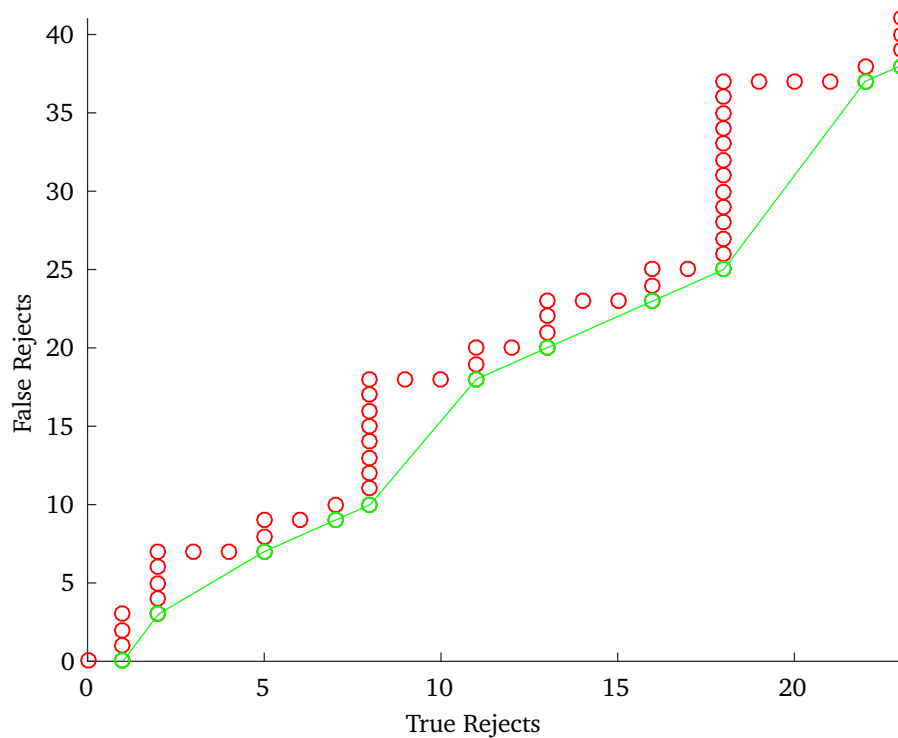
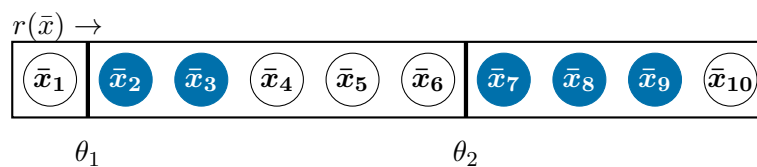
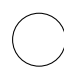



Fig. 2.4: This is an example of classified points sorted by a measure of confidence. You can see that $r(\bar{x}_2)$ and $r(\bar{x}_7)$ are the only sensible thresholds. Choosing $r(\bar{x}_8)$ instead of $r(\bar{x}_7)$ would only add an additional false reject and choosing $r(\bar{x}_6)$ over $r(\bar{x}_7)$ would result in one less true reject.



-  falsely classified points
-  correctly classified points

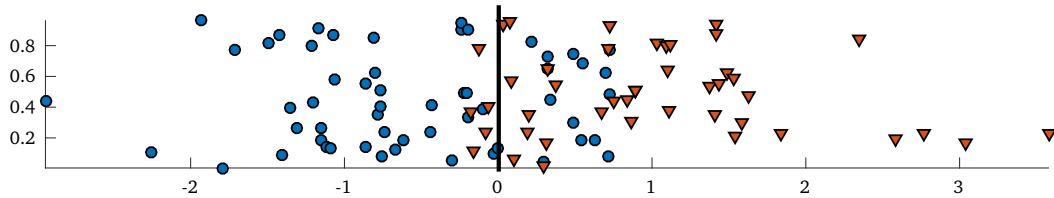
2.2.1 Global Reject

To adapt our reject strategy from before, we search for a threshold θ and reject according to our measures of confidence:

$$\max_i r_i(\bar{x}) < \theta : \bar{x} \text{ rejected}$$

We look again to choose θ so that it meets our requirements for an optimal threshold (see chapter 2.1.2). This strategy comes with the problem that it relies on all measures r_i to be scaled the same way (TODO: exmaple?). While there might be a workaround for this a global reject is still imprecise when the internal structures of the classes differ a lot(see figure 2.5). A reasonable threshold for a class with most points and also most classification errors near the decision plane is probably not well suited for a class with a wider (TODO: wording?) set of data points. This leads to the idea of having individual thresholds for each class.

Fig. 2.5: In this two-class classification example one class (blue circles) has few errors near the decision boundary and the other (red triangles) multiple farer away. A global reject optimal for the first class would result in a lot of unrejected errors in class 2. Conversely an optimal global reject for class 2 would result in many false rejects in class 1.



2.2.2 Local Reject

To account for differences in class structure each one now has a local reject threshold θ_i . A point is rejected if

$$\max_i r_i(\bar{x}) < \theta_i \mid \forall \bar{x} \text{ where } f(\bar{x}) = i$$

(TODO: need max?) This gives us an i-dimensional threshold vector $\bar{\theta} = (\theta_1, \dots, \theta_i)$. Each θ_i regulates the reject practice only in their respective class. (TODO: example)

2.2.3 Optimal Local Reject

For each threshold θ_i the same criteria apply to determine if it is optimal as for a global threshold (see chapter 2.1.2). With the difference that we look at each class individually. So the true and false rejects are now given by

$$T_{\theta_i}^i = R_{\theta_i}^i \cap E_i \quad (2.9)$$

$$F_{\theta_i}^i = R_{\theta_i}^i \cap L_i \quad (2.10)$$

where $R_{\theta_i}^i$ is the amount of rejected points in class i given the threshold θ_i , E_i the amount of falsely classified points in class i and L_i the amount of correctly classified points in i . As before we consider a certain set Θ_i of possible thresholds for each class (see chapter 2.1.3). We conclude that the threshold vector $\bar{\theta} = (\theta_1, \dots, \theta_N)$ is optimal if each θ_i is optimal.

2.2.4 Computation by Brute Force

To now find the optimal local reject vectors $\bar{\theta}$ with a brute force approach we need to consider every combination of thresholds in the sets Θ_i . Let

$$\mathbb{P} = \left\{ \bar{\theta} = (\theta_1, \dots, \theta_N) \in \Theta_1 \times \dots \times \Theta_N \right\}$$

be the set of all permutations of the sets of optimal thresholds from each class. Algorithm 1 describes in pseudo code how to find all $\bar{\theta}$.

Data: \mathbb{P}

Result: set of optimal local reject vectors Θ_{opt}

initialization;

array opt where $opt(n)$ is the maximum number of true rejects with n false ones ;

array $theta$ where $theta(n)$ is the optimal threshold vector with n false rejects ;

foreach $\bar{\theta} \in \mathbb{P}$ **do**

$F :=$ amount of false rejects for $\bar{\theta}$;

$T :=$ amount of true rejects for $\bar{\theta}$;

if $opt(F) < T$ **then**

$opt(F) := T$;

$theta(F) = \bar{\theta}$;

end

end

return $theta$

Algorithm 1: Computing optimal local reject options by brute force.

2.2.5 Computation by Dynamic Programming

Computation by brute force scales exponentially with the number of classes and is therefore not a feasible solution for large data sets with lots of clusters. But our problem is equivalent to a multiple choice knapsack problem (MCKP) where thresholds within a class correspond to items, false rejects correspond to costs, and true rejects correspond to their value (TODO:link). This allows us a faster solution that still maintains optimal results. Let

$$opt(n, j, i) = \max_{\bar{\theta}} \left\{ |T_{\bar{\theta}}| \text{ s.t. } \begin{array}{ll} |F_{\bar{\theta}}| = & n \\ \theta_k \in & \Theta_k \forall k < j \\ \theta_j \in & \{\theta_j(0), \dots, \theta_j(i)\} \\ \theta_k \in & \theta_k(0) \forall k > j \end{array} \right\}$$

(TODO: define $|T_{\bar{\theta}}|$ and $|F_{\bar{\theta}}|$???)

be the maximum number of true rejects with n false rejects while considering all thresholds in classes before class j and the first i thresholds in class j . All $\bar{\theta}$ corresponding to $opt(n, N, |\Theta_N| - 1) \forall n < |L|$ fulfill our criteria for an optimal threshold vector.

This results in the Bellmann equation 2.11 to efficiently compute opt :

$$opt(n, j, i) = \begin{cases} \text{if } n = 0 : & \sum_{k=1}^N |T_{\Theta_k(0)}^k| & (2.11a) \\ \text{if } n > 0, j = 0 : & \sum_{k=1}^N |T_{\Theta_k(0)}^k| & (2.11b) \\ \text{if } n > 0, j > 0, i = 0 : & opt(n, j-1, |\Theta_{j-1}| - 1) & (2.11c) \\ \text{if } 0 < n < |F_{\Theta_j(i)}^j|, j > 0 : & opt(n, j, i-1) & (2.11d) \\ \text{if } n \geq |F_{\Theta_j(i)}^j| > 0, j > 0 : & \max \left\{ \begin{array}{l} opt(n, j, i-1), \\ opt(n - |F_{\Theta_j(i)}^j|, j-1, |\Theta_{j-1}| - 1) + \\ |F_{\Theta_j(i)}^j| - |F_{\Theta_j(0)}^j| \end{array} \right\} & (2.11e) \end{cases}$$

Each case is explained as follows:

- Case 2.11a: $n = 0$ means that no false rejects are allowed so the first threshold (index 0) in each class is chosen. The resulting amount of true rejects is summed up over all classes.

- Case 2.11b: $j = 0$ means that no class is under consideration for a threshold. So we stay with the first threshold in each class (see above).
- Case 2.11c: $i = 0$ means that no threshold is under consideration in class j . So we look in the previous cell with the strictest threshold possible.
- Case 2.11d: The chosen threshold i in class j exceeds the allowed amount of false rejects, so the next less strict threshold is considered.
- Case 2.11e: Here the i th threshold in j is a possible threshold but it is not clear whether it is optimal. We consider both cases. If it is not the optimal threshold, we take the next less strict one. If it is optimal, we continue our search in the previous class but with $|F_{\Theta_j(i)}^j|$ less allowed false rejects in consequence to choosing this threshold. The other consequence is that this threshold results in a number of gained true rejects compared to the least strict threshold and this gain is added.

TODO: example (table?) compare to method in previous paper (loop count)

2.2.6 Greedy Computation

Computation by dynamic programming gives us an optimal local reject option, but it still might be unfeasible in some cases since it "scales quadratically with the number of data" (TODO: quote reject paper). Hence we are looking for a greedy approximation with linear running time. For this approach we define

$$\Delta T_{\theta_i(j)}^i = T_{\theta_i(j)}^i - T_{\theta_i(j-1)}^i \quad (2.12)$$

$$\Delta F_{\theta_i(j)}^i = F_{\theta_i(j)}^i - F_{\theta_i(j-1)}^i \quad (2.13)$$

as the amount of true and false rejects gained by a threshold compared to its predecessor and

$$g = \Delta T_{\theta_i}^i - \Delta F_{\theta_i}^i$$

as the local gain. The greedy approach lies within choosing the thresholds with the best local gain. Initially all thresholds are set to the most tolerant in each class. Looking at the next strict one in each class we pick the one with the highest local gain until the strictest possible thresholds are reached. The pair of true and false rejects $(T_{\theta_i}^i, F_{\theta_i}^i)$ is saved in each step if it is an improvement to existing solutions. Algorithm 2 details this procedure in pseudo code.

Data: sets of thresholds Θ_i
Result: set of optimal local reject vectors Θ_{opt}
initialization;
array opt where $opt(n)$ is the maximum number of true rejects with n false ones ;
array $theta$ where $theta(n)$ is the optimal threshold vector with n false rejects ;
 $\bar{\theta} = \{\theta_0(0), \dots, \theta_N(0)\}$;
 $F :=$ amount of false rejects for $\bar{\theta}$;
 $T :=$ amount of true rejects for $\bar{\theta}$;
while a more strict threshold in at least one class exists **do**
 foreach $\theta_i(j) \in \bar{\theta}$ **do**
 if $\exists \theta_i(j+1)$ **then**
 $g(i) = \Delta T_{\theta_i(j+1)}^i - \Delta F_{\theta_i(j+1)}^i$;
 else
 $g(i) = -\infty$;
 end
 end
 $\hat{i} = \arg \max_i g(i)$;
 set threshold $\theta_i(j)$ in $\bar{\theta}$ to $\theta_i(j+1)$;
 $T = T + \Delta T_{\theta_i(j+1)}^i$;
 $F = F + \Delta F_{\theta_i(j+1)}^i$;
 if $opt(F) < T$ **then**
 $opt(F) := T$;
 $theta(F) = \bar{\theta}$;
 end
end
return $theta$

Algorithm 2: Computing optimal local reject options by greedy evaluation.

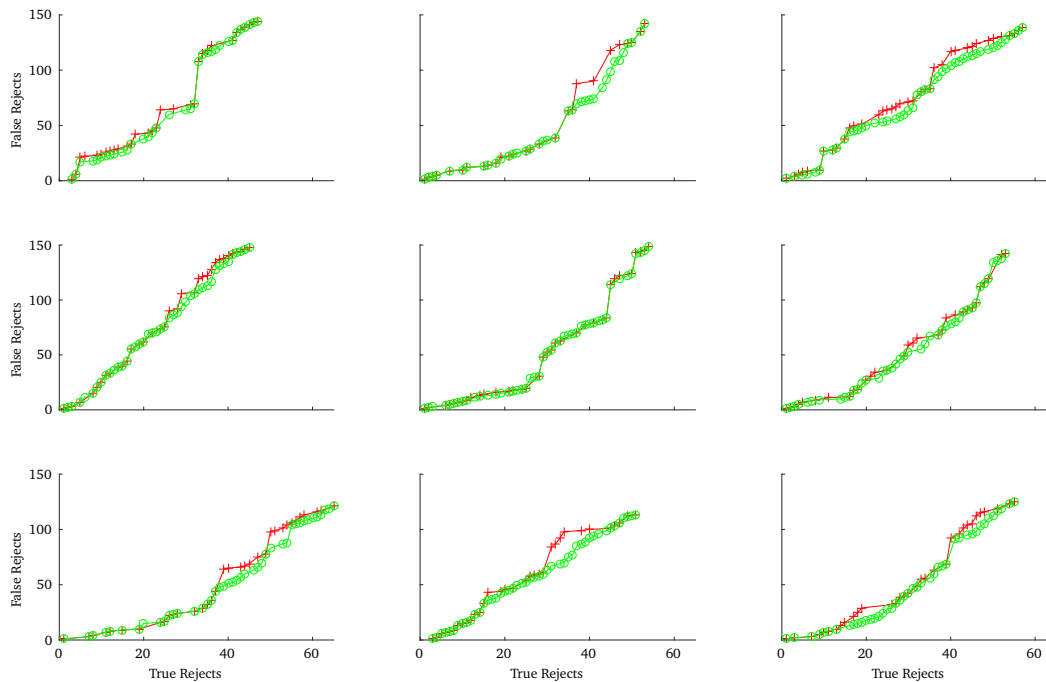
2.2.7 Evaluation

In this chapter we evaluate the described methods. We try to confirm the optimality of the dynamic programming method by comparing its result to the brute force algorithm. Furthermore we want to find out how close the greedy approach is to being optimal by comparing it to dynamic programming results.

DP vs Brute Force

DP vs Greedy

Fig. 2.6: Comparison of the thresholds computed greedily (red crosses) and the optimal ones computed by dynamic programming (green circles) on randomly generated data sets. We can see that the greedy strategy is close to optimal and mainly falls off slightly if a lot of points are rejected.



TODO:ARC

Application for SVM

3

Conclusions

4

List of Figures

2.1	Example of two classes separated by a decision boundary. You can see falsely classified points near the boundary which we would like to reject.	3
2.2	In this example we can see the corresponding true and false rejects if $r(\bar{x})$ is chosen as a threshold for each $\bar{x} \in X$. $r(\bar{x})$ is increasing from the bottom left corner to the upper right. As we can see, the number of false and true rejects are monotonically increasing with a more strict threshold.	5
2.3	This figure shows the thresholds from before (see figure 2.2). The green marked points correspond to threshold we consider optimal.	6
2.4	This is an example of classified points sorted by a measure of confidence. You can see that $r(\bar{x}_2)$ and $r(\bar{x}_7)$ are the only sensible thresholds. Choosing $r(\bar{x}_8)$ instead of $r(\bar{x}_7)$ would only add an additional false reject and choosing $r(\bar{x}_6)$ over $r(\bar{x}_7)$ would result in one less true reject.	6
2.5	In this two-class classification example one class (blue circles) has few errors near the decision boundary and the other (red triangles) multiple farer away. A global reject optimal for the first class would result in a lot of unrejected errors in class 2. Conversely an optimal global reject for class 2 would result in many false rejects in class 1.	7
2.6	Comparison of the thresholds computed greedily (red crosses) and the optimal ones computed by dynamic programming (green circles) on randomly generated data sets. We can see that the greedy strategy is close to optimal and mainly falls off slightly if a lot of points are rejected.	12

List of Tables

