

Maximum Margin Classification

Farzaneh Mirzazadeh

Thanks to Prof Dale Schuurmans

University of California, Santa Cruz

Winter '17

Maximum Margin Classification for Data with **Separable Classes**

Linearly separable classes

Assume data is linearly separable.
Consider entire set of consistent linear classifiers.

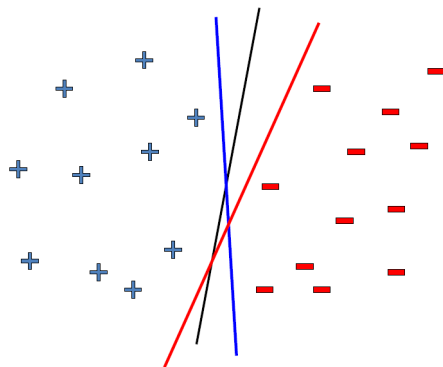
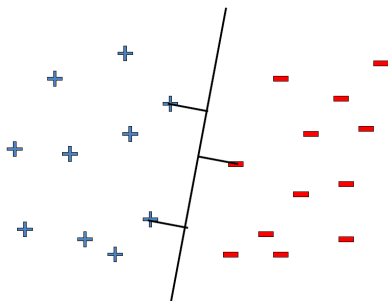


Figure from Prof Aarti Singh slides, CMU

Are some consistent linear classifiers better than others?
Which one we pick?

Linearly separable classes

- A simple but effective idea: **Maximum Margi linear classifier**
Pick the one with the largest margin!



8 Figure from Prof Aarti Singh slides, CMU

- Valadimir Vapnik (1970s in Soviet Union), brought the idea to US in 1990.
- Choose w, b to maximize the minimum distance between data points and decision boundary.

How can this be computed efficiently?

Six steps to show how:

- 1 Parameterizing (writing the equation for) the decision boundary
- 2 Writing mathematical meaning of correct classification for a point
- 3 Measuring confidence of classification (aka functional margin)
- 4 Measuring the geometric margin
- 5 Maximizing the margin
- 6 Re-express it as a quadratic program (QP)

Step 1: Parameterizing the decision boundary

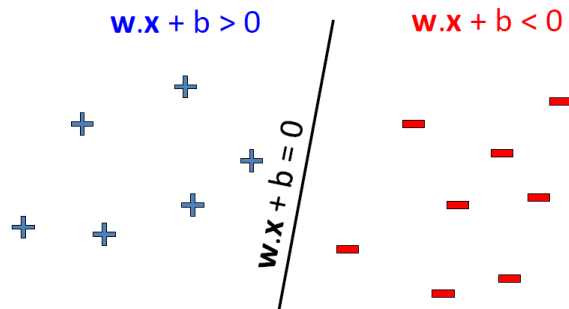


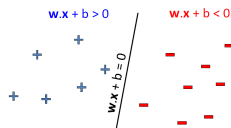
Figure from Prof Aarti Singh slides, CMU

- Equation of hyperplane: $\{\mathbf{x} : \mathbf{x}^\top \mathbf{w} + b = 0\}$

Note: Dot product= inner product $= \mathbf{x} \cdot \mathbf{w} = \mathbf{w} \cdot \mathbf{x} := \mathbf{x}^\top \mathbf{w} = \sum_{i=1}^n \mathbf{x}_i \mathbf{w}_i$

- A hyperplane separates the space to two half spaces.
 - Positive side $\mathbf{x}^\top \mathbf{w} + b > 0$
 - Negative side $\mathbf{x}^\top \mathbf{w} + b < 0$

Step 2. Correct classification



For any point $\langle \mathbf{x}_i, y_i \rangle$, $y_i \in \{-1, +1\}$
the hyperplane $\mathbf{x}^\top \mathbf{w} + b$ classifies it correctly means:

if $y_i > 0$ then $\mathbf{x}_i^\top \mathbf{w} + b > 0$

if $y_i < 0$ then $\mathbf{x}_i^\top \mathbf{w} + b < 0$

In summary

Correct classification for all training points $i = 1, 2, \dots, t$ requires
 $y_i(\mathbf{x}_i^\top \mathbf{w} + b) > 0, \quad \forall i.$

Step 3. Measuring confidence

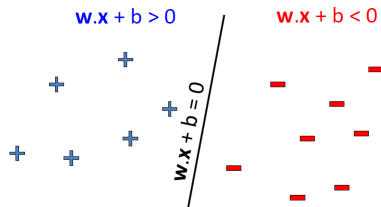


Figure from Prof Aarti Singh slides, CMU

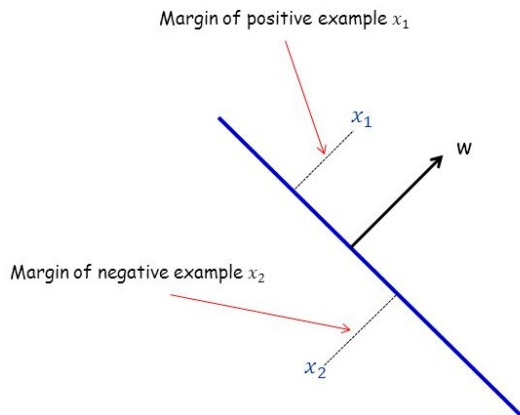
- **Question:** For a fixed hyperplane $\mathbf{x}^\top \mathbf{w} + b = 0$ what does the larger value of $(\mathbf{x}_i^\top \mathbf{w} + b)y_i$ show?
- **Answer:** Confidence of the classifier.
- **'Confidence'** = $(\mathbf{x}_i^\top \mathbf{w} + b)y_i$.
- Call this confidence, **functional margin**.

Step 4. Computing the margin (geometric)

Confidence: **functional margin** $= y(\mathbf{x}^\top \mathbf{w} + b)$ depends on scaling

geometric margin = distance between point and hyperplane

$$= \frac{y(\mathbf{x}^\top \mathbf{w} + b)}{\|\mathbf{w}\|}$$



Step 5. Maximizing the margin

Problem

Given a set of labeled examples, $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_t, y_t)$ where each $\mathbf{x}_i \in \mathbb{R}^n$ and each $y_i \in \{+1, -1\}$, find a weight vector \mathbf{w} and intercept b such that $\text{sign}(\mathbf{x}_i^\top \mathbf{w} + b) = y_i$ for all i .

Assume linearly separable! Want to maximize the minimum *margin* (over training point).

Maximum margin training

$$\max_{\mathbf{w}, b} \min_{i=1}^t \frac{y_i(\mathbf{x}_i^\top \mathbf{w} + b)}{\|\mathbf{w}\|}$$

How?

Recall the trick you have seen many times so far to convert min to a linear form. Use that. Then form a quadratic program (quadratic objective, linear constraints) with a little tweak. See next slides.

Step 6. Forming the quadratic program

$$\max_{\mathbf{w}, b} \min_{i=1}^t \text{ (geometric) margin of training point } i$$

$$\equiv \max_{\mathbf{w}, b} \min_{i=1}^t \frac{y_i(\mathbf{x}^\top \mathbf{w}_i + b)}{\|\mathbf{w}\|}$$

$$\equiv \max_{\delta, \mathbf{w}, b} \delta \quad \text{subject to} \quad \frac{y_i(\mathbf{x}^\top \mathbf{w}_i + b)}{\|\mathbf{w}\|} \geq \delta, \quad i = 1, 2, \dots, t$$

$$\equiv \max_{\delta, \mathbf{w}, b} \delta \quad \text{subject to} \quad y_i(\mathbf{x}^\top \mathbf{w}_i + b) \geq \delta \|\mathbf{w}\|, \quad i = 1, 2, \dots, t$$

$$\equiv \dots$$

If data linearly separable, then there exists \mathbf{w}, b s.t. $y_i(\mathbf{x}_i^\top \mathbf{w} - b) > 0, \forall i$. Thus any δ can be achieved by rescaling \mathbf{w}, b . So just rescale RHS to 1.

See next slide.

Step 6. Cont'd

...

$$\equiv \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \quad \text{subject to} \quad y_i(\mathbf{x}^\top \mathbf{w}_i + b) \geq 1, \quad i = 1, 2, \dots, t$$

$$\equiv \min_{\mathbf{w}, b} \|\mathbf{w}\| \quad \text{subject to} \quad y_i(\mathbf{x}^\top \mathbf{w}_i + b) \geq 1, \quad i = 1, 2, \dots, t$$

$$\equiv \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{x}^\top \mathbf{w}_i + b) \geq 1, \quad i = 1, 2, \dots, t$$

Quadratic objective with t linear constraints. **A Quadratic Program**

Maximum Margin Classification Training Optimization

Linearly separable case

Quadratic Programming

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \mathbf{w} + b) \geq 1, \quad i = 1, 2, \dots, t$$

Matrix form

All t constraint combined.

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w} \quad \text{s.t.} \quad \Delta(\mathbf{y})(X\mathbf{w} + \mathbf{1}b) \geq \mathbf{1}$$

where $\Delta(\cdot)$ is an operator that puts a vector on the diagonal of a diagonal matrix.
("diag" function in Matlab.)

Prediction

For a new unseen test point \mathbf{x}_o , predict as:

$$\hat{y}(\mathbf{x}_o) = \begin{cases} +1 & \mathbf{x}_o^\top \mathbf{w} \geq 0 \\ -1, & \mathbf{x}_o^\top \mathbf{w} < 0 \end{cases}$$

Support Vectors

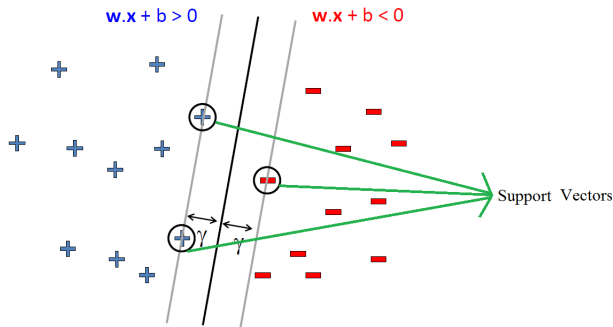


Figure from Prof Aarti Singh slides, CMU

Support Vectors: The subset of training examples that are **closest** to the optimal hyperplane defined by w , b .

This gives critical constraints! Removing other data doesn't change solution!

Support Vector Machines

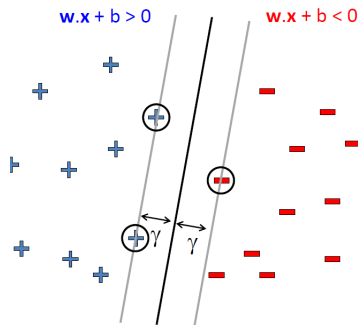


Figure from Prof Aarti Singh slides, CMU

- Data is summarized to its support vectors.
- Get a form of data compression.
- Provably good generalization if expected number of support vectors is much smaller than number of training examples.
- Thus called **Support Vector Machines** (SVMs).

General Case: Maximum Margin Classification for Non-separable Class Data

Non-separable case

- Noise and/or inconsistency in data
- Question: What if there is no consistent linear classifier
- Answer: Previous QP becomes infeasible. i.e will not have any solution
- Not good
- So what can one do?

Standard (heuristic) approach:

Add slack to margin

- Add a **slack variable** to each margin constraint.
- Try to maximize **minimum margin** minus a **sum of slacks**.
- Allow the maximum margin algorithm to **give up** on some bad training examples.

Much more robust against **class noise**!

Add slack to margins

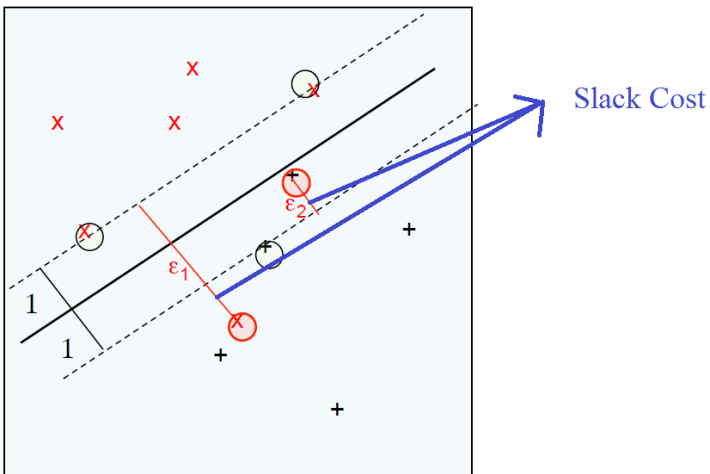


Figure from Prof David Helmbold Slides.

Quadratic Programming Algorithm with Slacks

$$\min_{\mathbf{w}, \mathbf{b}, \boldsymbol{\epsilon}} \quad \frac{\beta}{2} \|\mathbf{w}\|_2^2 + \mathbf{1}^\top \boldsymbol{\epsilon} \quad \text{s.t.} \quad \boldsymbol{\epsilon} \geq \mathbf{1} - \Delta(\mathbf{y})(X\mathbf{w} + \mathbf{1}b) \\ \boldsymbol{\epsilon} \geq \mathbf{0}$$

- Add a slack variable $\epsilon_i \geq 0$ for each training example (X_i, y_i)
- Allow slack to reduce the margin, but at a cost (linear in mistake)
- β regularization (trade-off) parameter
- Called **Soft margin Support Vector Machines**

Hinge loss

- Note: The quadratic program that you get is exactly equivalent to minimizing a L_2 regularized form + a loss function

$$\min_{\mathbf{w}, b} \frac{\beta}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^t \max(0, 1 - y_i(\mathbf{X}_i \mathbf{w} + b))$$

- The loss is called **Hinge Loss**

$$\text{hinge}(m) = \begin{cases} 1 - m & m \leq 1 \\ 0 & m > 1 \end{cases}$$

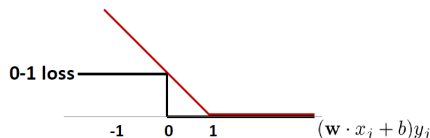


Figure from Prof Aarti Singh slides, CMU

- Again a loss + regularizer minimization as before

Summary

- We learned maximum margin classification:
 - Case 1: Linearly separable (hard margin SVM)
 - Case 2: Linearly non-separable (soft margin SVM)
- Showed that soft margin SVM is a regularized hinge loss.

Where to study?

- James et al (2013) Chap 9
- Hastie et al. 2nd ed (2009) Sec4.5.2, 12.2, 12.3
- Bishop (2006) Sec 7.1
- Duda et al. 2nd ed(2001) Chap 5
- Cherkassky & Mulier (1998) Chap 9
- Marsland (2009) Chap 5