

Learning Maximum-Margin Hyperplanes: Support Vector Machines

hard margin

Piyush Rai

soft margin

hands on
session

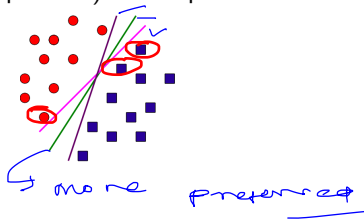
Machine Learning (CS771A)

Aug 24, 2016

Perceptron and (Lack of) Margins

- Perceptron learns a hyperplane (of many possible) that separates the classes

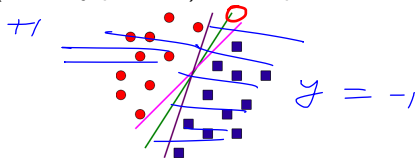
2 - class



green — max degree of
sep

Perceptron and (Lack of) Margins

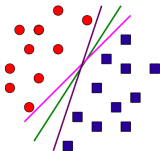
- Perceptron learns a hyperplane (of many possible) that separates the classes



- Standard Perceptron doesn't guarantee any margin around the hyperplane

Perceptron and (Lack of) Margins

- Perceptron learns a hyperplane (of many possible) that separates the classes



- Standard Perceptron doesn't guarantee any “margin” around the hyperplane
- Note: Possible to “artificially” introduce a margin in the Perceptron

$$\vec{a} = [1, 2]$$

$$\vec{b} = [2, 3]$$

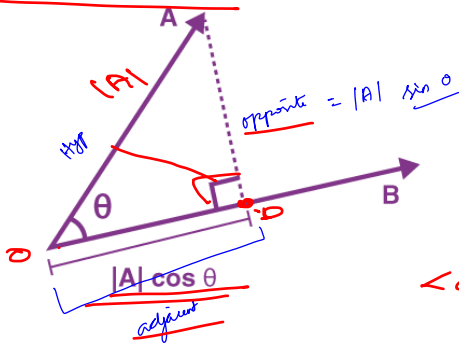
$$\vec{a} \cdot \vec{b} = \sum_i a_i b_i = 2 + 6 = 8$$

$$= \|\vec{a}\| \|\vec{b}\| \cos \theta$$

dot products and margins

DOT PRODUCT OF VECTORS

BYJU'S
The Learning App



$$\sin \theta = \frac{\text{opp}}{\text{hyp}} = \frac{BP}{OA}$$

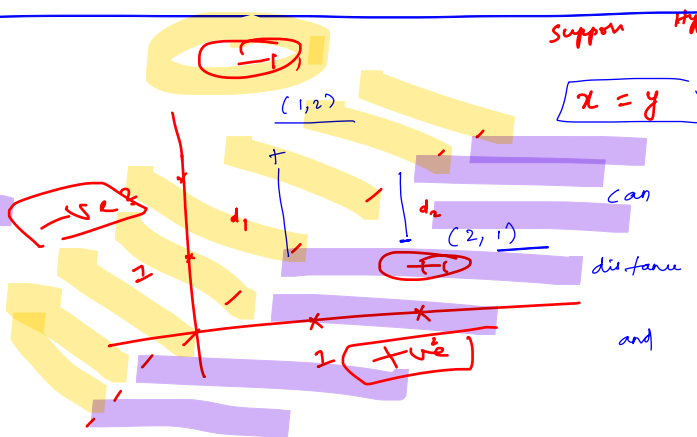
$$= \frac{BP}{OA}$$

$$\cos \theta = \frac{OP}{OA}$$

$$\vec{a} \cdot \vec{b} = |a| |b| \cos \theta$$

Semantics of

Signs (wrt x)



Support hyper plane is

$$x = y$$

can we compute the distance b/w point and the hyper plane

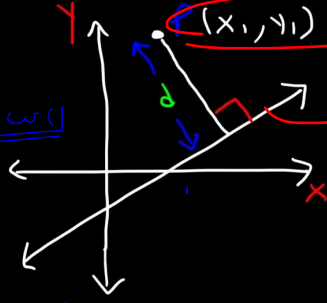
Point - Line Distance

$$d = \frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

$\|w\|$

$$ax + by + c = 0$$

$$P(x_1, y_1)$$



$$ax + by + c$$

Hyper plane is

$$x - y = 0$$

$$a = 1 \quad \text{and} \quad b = -1 \quad \underline{c = 0}$$

d1

d2

$$d_1 = \frac{\cancel{1} - 2 + 0 \cancel{0}}{\sqrt{1+4}} = \frac{-1}{\sqrt{5}}$$

$$d_2 = \frac{\cancel{2} - 1 + 0 \cancel{0}}{\sqrt{5}} = \frac{+1}{\sqrt{5}}$$

dot product b/w line and point

denote \vec{w} as coefficients of the line

$$\vec{w} = \begin{bmatrix} \overset{x}{1}, & \overset{y}{-1}, & \overset{c}{0} \\ a & b & c \end{bmatrix}$$

$$\begin{aligned} P_1 &= (1, 1) \\ P_2 &= (2, 1) \end{aligned}$$

augment each point with 1 towards the

end to account for 'c'

$$x_1 = \begin{bmatrix} 1, 2, \underline{1} \end{bmatrix} \quad \text{and} \quad x_2 = \begin{bmatrix} 2, 1, \underline{1} \end{bmatrix}$$

$$\underline{w^T x_1} = 1 - 2 + 0 = \frac{-1}{\sqrt{5}} \quad \text{and} \quad w^T x_2 = 2 - 1 + 0 = \frac{+1}{\sqrt{5}}$$

This is unnormalized signed distance.

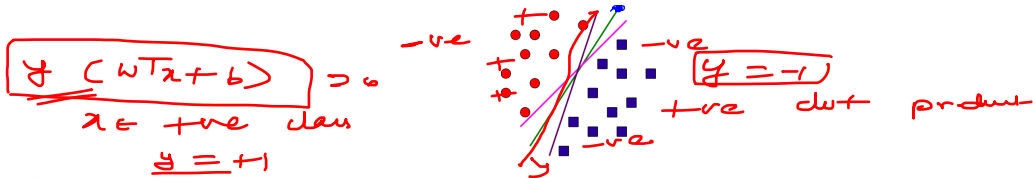
$(1, 2)$

$$\text{Norm factor} = \frac{1}{\|w\|} = \frac{1}{\sqrt{1^2 + 2^2 + 0^2}} = \frac{1}{\sqrt{5}}$$

Sign Indicates the side of the margin

Perceptron and (Lack of) Margins

- Perceptron learns a hyperplane (of many possible) that separates the classes



- Standard Perceptron doesn't guarantee any "margin" around the hyperplane
- Note: Possible to "artificially" introduce a margin in the Perceptron

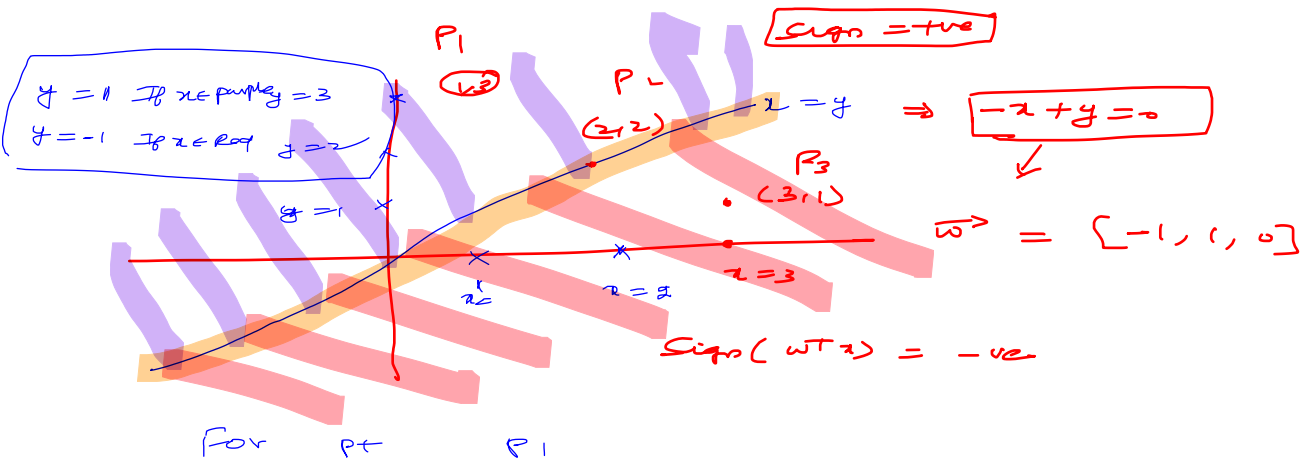
- Simply change the Perceptron mistake condition to

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq \gamma$$

$$y(\mathbf{w}^T \mathbf{x} + b) > \gamma$$

where $\gamma > 0$ is a **pre-specified margin**. For standard Perceptron, $\gamma = 0$

$\forall x \in \text{Pattern}$



$$w^T P_1 = -1 + 3 + 0 = +2$$

$$\text{For } P_2 \quad w^T P_2 = -2 + 2 = 0 \quad \left[w^T P_3 = -2 \right]$$

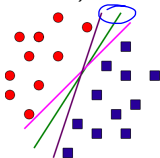
Solve for a \vec{w} such that

$$y (w^T x + b) > 0 \quad \forall x \in \text{Train data}$$

unnormalized distance

Perceptron and (Lack of) Margins

- Perceptron learns a hyperplane (of many possible) that separates the classes



→ → →

- Standard Perceptron doesn't guarantee any “margin” around the hyperplane
- Note: Possible to “artificially” introduce a margin in the Perceptron
 - Simply change the Perceptron mistake condition to

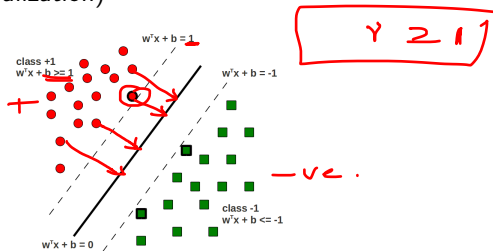
$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq \gamma$$

where $\gamma > 0$ is a **pre-specified margin**. For standard Perceptron, $\gamma = 0$

- **Support Vector Machine (SVM)** offers a more principled way of doing this by learning the **maximum margin hyperplane**

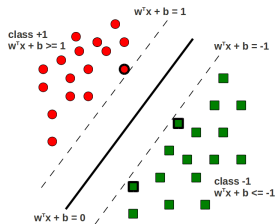
Support Vector Machine (SVM)

- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



Support Vector Machine (SVM)

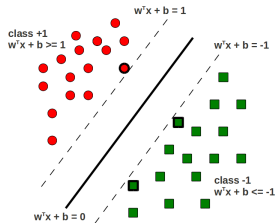
- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



- SVMs can also learn **nonlinear decision boundaries** using **kernels** (though the idea of kernels is not specific to SVMs and is more generally applicable)

Support Vector Machine (SVM)

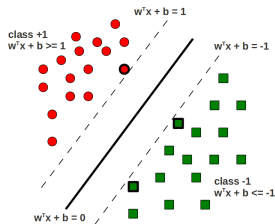
- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



- SVMs can also learn **nonlinear decision boundaries** using **kernels** (though the idea of kernels is not specific to SVMs and is more generally applicable)
- Reason behind the name “Support Vector Machine”?

Support Vector Machine (SVM)

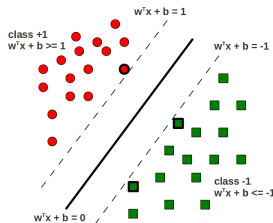
- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



- SVMs can also learn **nonlinear decision boundaries** using **kernels** (though the idea of kernels is not specific to SVMs and is more generally applicable)
- Reason behind the name “Support Vector Machine”? SVM finds the most important examples (called “**support vectors**”) in the training data

Support Vector Machine (SVM)

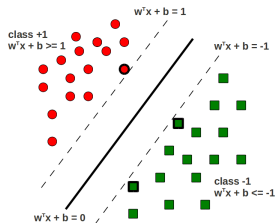
- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



- SVMs can also learn **nonlinear decision boundaries** using **kernels** (though the idea of kernels is not specific to SVMs and is more generally applicable)
- Reason behind the name “Support Vector Machine”? SVM finds the most important examples (called **“support vectors”**) in the training data
 - These examples also “balance” the margin boundaries (hence called “support”).

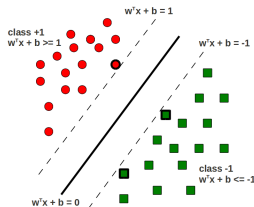
Support Vector Machine (SVM)

- Learns a hyperplane such that the positive and negative class training examples are **as far away as possible** from it (ensures good generalization)



- SVMs can also learn **nonlinear decision boundaries** using **kernels** (though the idea of kernels is not specific to SVMs and is more generally applicable)
- Reason behind the name “Support Vector Machine”? SVM finds the most important examples (called **“support vectors”**) in the training data
 - These examples also “balance” the margin boundaries (hence called “support”). Also, even if we throw away the remaining training data and re-learn the SVM classifier, we’ll get the same hyperplane

Learning a Maximum Margin Hyperplane



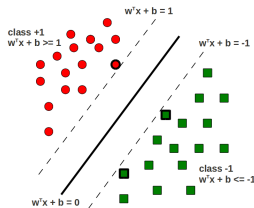
$$(\vec{w}, b)$$

- Suppose there exists a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ such that

- $\mathbf{w}^T \mathbf{x}_n + b \geq 1$ for $y_n = +1$
- $\mathbf{w}^T \mathbf{x}_n + b \leq -1$ for $y_n = -1$

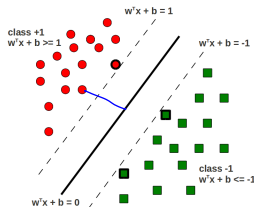
$$y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$$

Learning a Maximum Margin Hyperplane



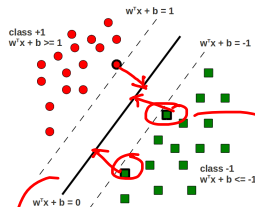
- Suppose there exists a hyperplane $w^T x + b = 0$ such that
 - $w^T x_n + b \geq 1$ for $y_n = +1$
 - $w^T x_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(w^T x_n + b) \geq 1 \quad \forall n$ (the margin condition)

Learning a Maximum Margin Hyperplane



- Suppose there exists a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ such that
 - $\mathbf{w}^T \mathbf{x}_n + b \geq 1$ for $y_n = +1$
 - $\mathbf{w}^T \mathbf{x}_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n$ (the margin condition)
 - Also note that $\min_{1 \leq n \leq N} |\mathbf{w}^T \mathbf{x}_n + b| = 1$

Learning a Maximum Margin Hyperplane



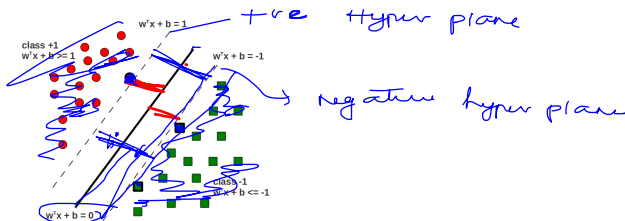
$$|w^T x + b| = \text{distance}$$

- Suppose there exists a hyperplane $w^T x + b = 0$ such that
 - $w^T x_n + b \geq 1$ for $y_n = +1$
 - $w^T x_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(w^T x_n + b) \geq 1 \quad \forall n$ (the margin condition)
 - Also note that $\min_{1 \leq n \leq N} |w^T x_n + b| = 1$
 - Thus margin on each side: $\gamma = \min_{1 \leq n \leq N} \frac{|w^T x_n + b|}{\|w\|}$

$$\max \frac{1}{\|w\|}$$

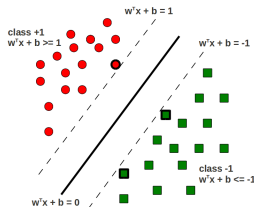
$$\text{st } y_n(w^T x_n + b) \geq 1 \quad \forall n$$

Learning a Maximum Margin Hyperplane



- Suppose there exists a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ such that
 - $\mathbf{w}^T \mathbf{x}_n + b \geq 1$ for $y_n = +1$
 - $\mathbf{w}^T \mathbf{x}_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n$ (the margin condition)
 - Also note that $\min_{1 \leq n \leq N} |\mathbf{w}^T \mathbf{x}_n + b| = 1$
 - Thus margin on each side: $\gamma = \min_{1 \leq n \leq N} \frac{|\mathbf{w}^T \mathbf{x}_n + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$
 - Total margin $= 2\gamma = \frac{2}{\|\mathbf{w}\|}$

Learning a Maximum Margin Hyperplane



- Suppose there exists a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ such that
 - $\mathbf{w}^T \mathbf{x}_n + b \geq 1$ for $y_n = +1$
 - $\mathbf{w}^T \mathbf{x}_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n$ (the margin condition)
 - Also note that $\min_{1 \leq n \leq N} |\mathbf{w}^T \mathbf{x}_n + b| = 1$
 - Thus margin on each side: $\gamma = \min_{1 \leq n \leq N} \frac{|\mathbf{w}^T \mathbf{x}_n + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$
 - Total margin = $2\gamma = \frac{2}{\|\mathbf{w}\|}$
- Want the hyperplane (\mathbf{w}, b) to have the largest possible margin

Large Margin = Good Generalization

- Large margins intuitively mean good generalization

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|w\|}$

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}
- Small $\|\mathbf{w}\| \Rightarrow$ regularized/simple solutions (w_i 's don't become too large)

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}
- Small $\|\mathbf{w}\| \Rightarrow$ regularized/simple solutions (w_i 's don't become too large)
 - Recall our discussion of regularization..

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}
- Small $\|\mathbf{w}\| \Rightarrow$ regularized/simple solutions (w_i 's don't become too large)
 - ~~Recall our discussion of regularization..~~
- Simple solutions \Rightarrow good generalization on test data

Large Margin = Good Generalization

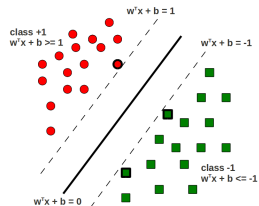
- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}
- Small $\|\mathbf{w}\| \Rightarrow$ regularized/simple solutions (w_i 's don't become too large)
 - Recall our discussion of regularization..
- Simple solutions \Rightarrow good generalization on test data
- Want to see an even more formal justification? :-)

Large Margin = Good Generalization

- Large margins intuitively mean good generalization
- We saw that margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
- Large margin \Rightarrow small $\|\mathbf{w}\|$, i.e., small ℓ_2 norm of \mathbf{w}
- Small $\|\mathbf{w}\| \Rightarrow$ regularized/simple solutions (w_i 's don't become too large)
 - Recall our discussion of regularization..
- Simple solutions \Rightarrow good generalization on test data
- Want to see an even more formal justification? :-)
 - Wait until we cover Learning Theory!

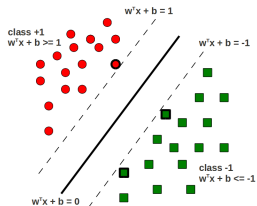
Hard-Margin SVM

- Every training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



Hard-Margin SVM

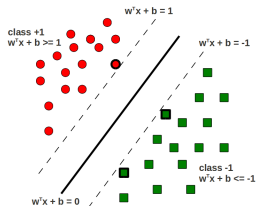
- Every training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



- Also want to maximize the margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$

Hard-Margin SVM

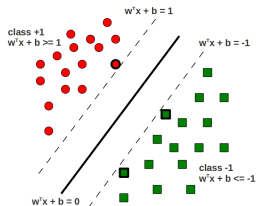
- Every training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



- Also want to maximize the margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
 - Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$

Hard-Margin SVM

- Every training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



- Also want to maximize the margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$

- Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$

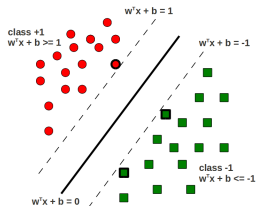
- The objective for hard-margin SVM

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \quad \rightarrow \text{quadratic cost function} \\ \text{subject to} \quad & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N \quad \rightarrow \text{linear constraint function} \end{aligned}$$

opt problem to be solved

Hard-Margin SVM

- Every training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



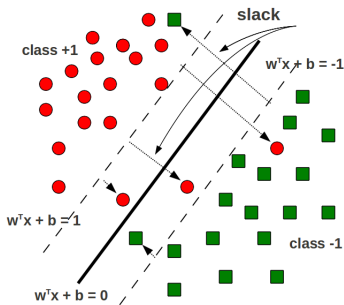
- Also want to maximize the margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$
 - Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$
- The objective for hard-margin SVM

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N \end{aligned}$$

- Thus the hard-margin SVM minimizes a **convex objective function** which is a **Quadratic Program (QP)** with N linear inequality constraints

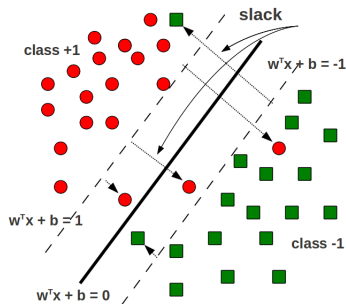
Soft-Margin SVM (More Commonly Used)

- Allow some training examples to fall **within the margin region**, or be even **misclassified** (i.e., fall on the wrong side). Preferable if **training data is noisy**



Soft-Margin SVM (More Commonly Used)

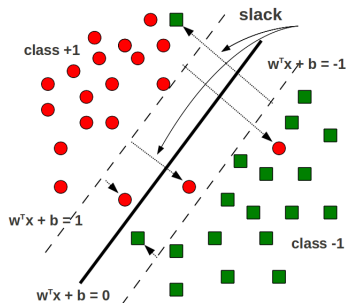
- Allow some training examples to fall **within the margin region**, or be even **misclassified** (i.e., fall on the wrong side). Preferable if **training data is noisy**



- Each training example (\mathbf{x}_n, y_n) given a “slack” $\xi_n \geq 0$ (distance by which it “violates” the margin). If $\xi_n > 1$ then \mathbf{x}_n is totally on the wrong side

Soft-Margin SVM (More Commonly Used)

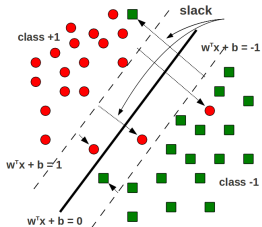
- Allow some training examples to fall **within the margin region**, or be even **misclassified** (i.e., fall on the wrong side). Preferable if **training data is noisy**



- Each training example (\mathbf{x}_n, y_n) given a “slack” $\xi_n \geq 0$ (distance by which it “violates” the margin). If $\xi_n > 1$ then \mathbf{x}_n is totally on the wrong side
 - Basically, we want a **soft-margin condition**: $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0$

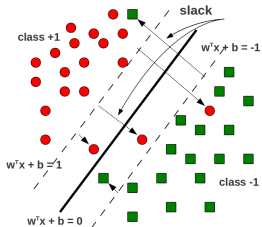
Soft-Margin SVM (More Commonly Used)

- Goal: Maximize the margin, while also minimizing the **sum of slacks** (don't want too many training examples violating the margin condition)



Soft-Margin SVM (More Commonly Used)

- Goal: Maximize the margin, while also minimizing the **sum of slacks** (don't want too many training examples violating the margin condition)



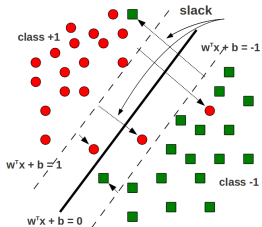
- The primal objective for soft-margin SVM can thus be written as

$$\min_{\mathbf{w}, b, \xi} f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n$$

subject to constraints $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N$

Soft-Margin SVM (More Commonly Used)

- Goal: Maximize the margin, while also minimizing the **sum of slacks** (don't want too many training examples violating the margin condition)



- The primal objective for soft-margin SVM can thus be written as

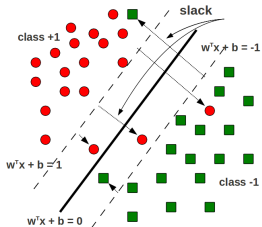
$$\min_{\mathbf{w}, b, \xi} f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n$$

subject to constraints $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N$

- Thus the soft-margin SVM also minimizes a **convex objective function** which is a **Quadratic Program** (QP) with $2N$ linear inequality constraints

Soft-Margin SVM (More Commonly Used)

- Goal: Maximize the margin, while also minimizing the **sum of slacks** (don't want too many training examples violating the margin condition)



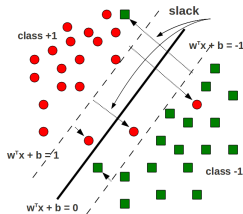
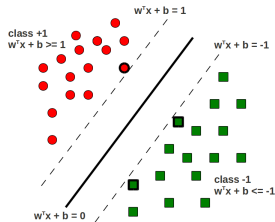
- The primal objective for soft-margin SVM can thus be written as

$$\min_{\mathbf{w}, b, \xi} f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n$$

subject to constraints $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N$

- Thus the soft-margin SVM also minimizes a **convex objective function** which is a **Quadratic Program** (QP) with $2N$ linear inequality constraints
- Param. C controls the trade-off between large margin vs small training error

Summary: Hard-Margin SVM vs Soft-Margin SVM



- Objective for the hard-margin SVM (unknowns are \mathbf{w} and b)

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to constraints} \quad & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N \end{aligned}$$

- Objective for the soft-margin SVM (unknowns are \mathbf{w} , b , and $\{\xi_n\}_{n=1}^N$)

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N \end{aligned}$$

- In either case, we have to solve constrained, convex optimization problem

Brief Detour: Solving Constrained Optimization Problems

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & \cancel{h_m(\mathbf{w}) = 0}, \quad m = \cancel{1}, \dots, \cancel{M} \end{aligned}$$

let us ignore the equality constraints since SVM
does not have them

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

- Consider $\mathcal{L}_P(\mathbf{w}) = \max_{\alpha, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$. Note that

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

- Consider $\mathcal{L}_P(\mathbf{w}) = \max_{\alpha, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$. Note that

- $\mathcal{L}_P(\mathbf{w}) = \infty$ if \mathbf{w} violates any of the constraints (g 's or h 's)

assume $g_n(\mathbf{w})$ is violated

$$\Rightarrow g_n(\mathbf{w}) > 0$$

Now when we max $\alpha_n g_n(\mathbf{w})$ we can trivially

set $\alpha_n = +\infty$ and the obj is maximized

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

- Consider $\mathcal{L}_P(\mathbf{w}) = \max_{\alpha, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$. Note that

- $\mathcal{L}_P(\mathbf{w}) = \infty$ if \mathbf{w} violates any of the constraints (g 's or h 's) *If $g_n(\mathbf{w})$ is satisfied*
- $\mathcal{L}_P(\mathbf{w}) = f(\mathbf{w})$ if \mathbf{w} satisfies all the constraints (g 's and h 's)

$g_n(\mathbf{w}) \leq 0$
what is α_n in $\max_{\alpha_n} \alpha_n g_n(\mathbf{w})$?

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

- Consider $\mathcal{L}_P(\mathbf{w}) = \max_{\alpha, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$. Note that
 - $\mathcal{L}_P(\mathbf{w}) = \infty$ if \mathbf{w} violates any of the constraints (g 's or h 's)
 - $\mathcal{L}_P(\mathbf{w}) = f(\mathbf{w})$ if \mathbf{w} satisfies all the constraints (g 's and h 's)
- Thus $\min_{\mathbf{w}} \mathcal{L}_P(\mathbf{w}) = \min_{\mathbf{w}} \max_{\alpha \geq 0, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$ solves the same problem as the original problem and will have the same solution. For convex f, g, h , the order of min and max is interchangeable.

Constrained Optimization via Lagrangian

- Consider optimizing the following objective, subject to some constraints

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & g_n(\mathbf{w}) \leq 0, \quad n = 1, \dots, N \\ & h_m(\mathbf{w}) = 0, \quad m = 1, \dots, M \end{aligned}$$

- Introduce Lagrange multipliers $\alpha = \{\alpha_n\}_{n=1}^N$, $\alpha_n \geq 0$, and $\beta = \{\beta_m\}_{m=1}^M$, one for each constraint, and construct the following Lagrangian

$$\mathcal{L}(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{n=1}^N \alpha_n g_n(\mathbf{w}) + \sum_{m=1}^M \beta_m h_m(\mathbf{w})$$

- Consider $\mathcal{L}_P(\mathbf{w}) = \max_{\alpha, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$. Note that
 - $\mathcal{L}_P(\mathbf{w}) = \infty$ if \mathbf{w} violates any of the constraints (g 's or h 's)
 - $\mathcal{L}_P(\mathbf{w}) = f(\mathbf{w})$ if \mathbf{w} satisfies all the constraints (g 's and h 's)
- Thus $\min_{\mathbf{w}} \mathcal{L}_P(\mathbf{w}) = \min_{\mathbf{w}} \max_{\alpha \geq 0, \beta} \mathcal{L}(\mathbf{w}, \alpha, \beta)$ solves the same problem as the original problem and will have the same solution. For convex f, g, h , the order of min and max is interchangeable.
- Karush-Kuhn-Tucker (KKT) Conditions:** At the optimal solution, $\alpha_n g_n(\mathbf{w}) = 0$ (note the \max_{α})

Solving Hard-Margin SVM

Solving Hard-Margin SVM

- The hard-margin SVM optimization problem is:

$$\begin{array}{ll} \min_{\mathbf{w}, b} & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} & 1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0, \quad n = 1, \dots, N \end{array}$$

- A constrained optimization problem. Can solve using Lagrange's method

Solving Hard-Margin SVM

- The hard-margin SVM optimization problem is:

$$\begin{array}{ll} \min_{\mathbf{w}, b} & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} & 1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0, \quad n = 1, \dots, N \end{array}$$

- A constrained optimization problem. Can solve using Lagrange's method
- Introduce **Lagrange Multipliers** α_n ($n = \{1, \dots, N\}$), one for each constraint, and solve the following Lagrangian:

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{\|\mathbf{w}\|^2}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

- Note: $\alpha = [\alpha_1, \dots, \alpha_N]$ is the vector of Lagrange multipliers

Solving Hard-Margin SVM

- The hard-margin SVM optimization problem is:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & 1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0, \quad n = 1, \dots, N \end{aligned}$$

- A constrained optimization problem. Can solve using Lagrange's method
- Introduce **Lagrange Multipliers** α_n ($n = \{1, \dots, N\}$), one for each constraint, and solve the following Lagrangian:

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{\|\mathbf{w}\|^2}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

- Note: $\alpha = [\alpha_1, \dots, \alpha_N]$ is the vector of Lagrange multipliers
- We will solve this Lagrangian by solving a **dual problem** (eliminate \mathbf{w} and b and solve for the “**dual variables**” α)

Solving Hard-Margin SVM

- The original Lagrangian is

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{\mathbf{w}^T \mathbf{w}}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

For convex function $\min_{\mathbf{w}, b} \max_{\alpha} = \max_{\alpha} \min_{\mathbf{w}}$

Interested students can further read on "strong duality".

Solving Hard-Margin SVM

- The original Lagrangian is

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{\mathbf{w}^\top \mathbf{w}}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^\top \mathbf{x}_n + b)\}$$

- Take (partial) derivatives of \mathcal{L} w.r.t. \mathbf{w} , b and set them to zero

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0$$



This means that \vec{w} = weighted sum
of all points

!!! to perception

Solving Hard-Margin SVM

- The original Lagrangian is

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{\mathbf{w}^T \mathbf{w}}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

- Take (partial) derivatives of \mathcal{L} w.r.t. \mathbf{w} , b and set them to zero

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0$$

- Important: Note the form of the solution \mathbf{w} - it is simply a **weighted sum of all the training inputs** $\mathbf{x}_1, \dots, \mathbf{x}_N$ (and α_n is like the “importance” of \mathbf{x}_n)
- Substituting $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ in Lagrangian and also using $\sum_{n=1}^N \alpha_n y_n = 0$

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n y_m y_n (\mathbf{x}_m^T \mathbf{x}_n) \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

[†] If interested in more details of the solver, see: "Support Vector Machine Solvers" by Bottou and Lin

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is **maximizing a concave function** (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original primal SVM objective was also convex

[†] If interested in more details of the solver, see: "Support Vector Machine Solvers" by Bottou and Lin

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is **maximizing a concave function** (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original primal SVM objective was also convex
- **Important:** Inputs \mathbf{x} 's only appear as **inner products** (helps to “kernelize”)

[†] If interested in more details of the solver, see: “Support Vector Machine Solvers” by Bottou and Lin

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\boxed{\begin{array}{l} \max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0 \end{array}}$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is maximizing a concave function (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original primal SVM objective was also convex
- **Important:** Inputs \mathbf{x} 's only appear as inner products (helps to “kernelize”)
- Can solve[†] the above objective function for α using various methods, e.g.,

[†] If interested in more details of the solver, see: “Support Vector Machine Solvers” by Bottou and Lin

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\boxed{\begin{array}{ll} \max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha & \text{s.t. } \sum_{n=1}^N \alpha_n y_n = 0 \end{array}}$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is **maximizing a concave function** (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original primal SVM objective was also convex
- **Important:** Inputs \mathbf{x} 's only appear as **inner products** (helps to “kernelize”)
- Can solve[†] the above objective function for α using various methods, e.g.,
 - Treating the objective as a **Quadratic Program** (QP) and running some off-the-shelf QP solver such as quadprog (MATLAB), CVXOPT, CPLEX, etc.

[†] If interested in more details of the solver, see: “Support Vector Machine Solvers” by Bottou and Lin

Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is **maximizing a concave function** (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original primal SVM objective was also convex
- **Important:** Inputs \mathbf{x} 's only appear as **inner products** (helps to “kernelize”)
- Can solve[†] the above objective function for α using various methods, e.g.,
 - Treating the objective as a **Quadratic Program** (QP) and running some off-the-shelf QP solver such as quadprog (MATLAB), CVXOPT, CPLEX, etc.
 - Using **(projected) gradient methods** (projection needed because the α 's are constrained). Gradient methods will usually be much faster than QP methods.

[†] If interested in more details of the solver, see: “Support Vector Machine Solvers” by Bottou and Lin

Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n:y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n:y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

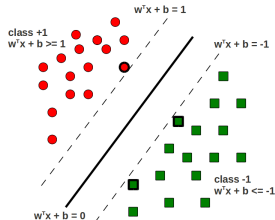
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n:y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n:y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**

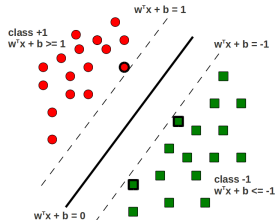
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n:y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n:y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**
- For the optimal α_n 's

$$\alpha_n \{1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$

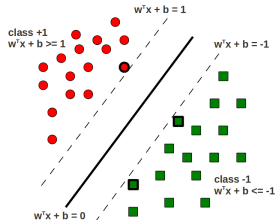
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n:y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n:y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT)** conditions
- For the optimal α_n 's
$$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$
- α_n is **non-zero** only if \mathbf{x}_n

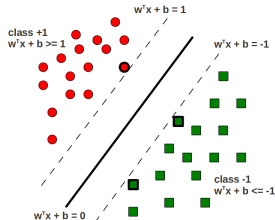
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n: y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n: y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**
- For the optimal α_n 's

$$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$

- α_n is **non-zero** only if \mathbf{x}_n lies on one of the two **margin boundaries**, i.e., for which $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$

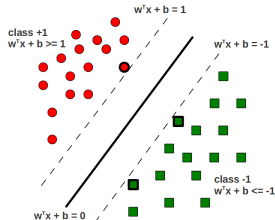
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n: y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n: y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**
- For the optimal α_n 's

$$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$

- α_n is **non-zero** only if \mathbf{x}_n lies on one of the two **margin boundaries**, i.e., for which $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
- These examples are called **support vectors**

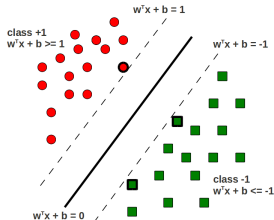
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

$$b = -\frac{1}{2} \left(\min_{n: y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n: y_n=-1} \mathbf{w}^T \mathbf{x}_n \right)$$

- A nice property:** Most α_n 's in the solution will be zero (**sparse solution**)



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**
- For the optimal α_n 's

$$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$

- α_n is **non-zero** only if \mathbf{x}_n lies on one of the two **margin boundaries**, i.e., for which $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
- These examples are called **support vectors**
- Recall the support vectors “support” the margin boundaries