

FSAN/ELEG815: Statistical Learning Gonzalo R. Arce

Department of Electrical and Computer Engineering University of Delaware

Support Vector Machines



Support Vector Machines - Better linear separation



Two questions:

- 1. Why is bigger margin better?
- 2. Which w maximizes the margin?



Support Vector Machines - Growth Function

All Possible Dichotomies with a line.



Bad news!



FSAN/ELEG815

Support Vector Machines - Growth Function



Let's consider a classifier that requires a minimum margin.



FSAN/ELEG815

Support Vector Machines - Growth Function



Let's consider a classifier that requires a minimum margin.

Fat margins imply fewer dichotomies \implies smaller growth function



Support Vector Machines - Finding ${f w}$ with large margin

Let \mathbf{x}_n be the nearest data point to the line/plane (given by $\mathbf{w}^\top \mathbf{x} = 0$)

How far is it?

Two preliminary techniques:

1. Normalize \mathbf{w} : For any point:

 $|\mathbf{w}^{\top}\mathbf{x}_n| > 0.$

Does scalar multiplication change the plane? NO! Pick one:

$$|\mathbf{w}^{\top}\mathbf{x}_n| = 1.$$

2. Pull out w_0 :

 $\mathbf{w} = (w_1, ... w_d)$ apart from $w_0 = b$. The plane is now $\mathbf{wx} + b = 0$ (no x_0).



The distance between \mathbf{x}_n and the plane $\mathbf{w}^\top \mathbf{x} + b = 0$, where $|\mathbf{w}^\top \mathbf{x}_n + b| = 1$.



The vector ${\bf w}$ is \perp to the plane in the ${\cal X}$ space:



The distance between \mathbf{x}_n and the plane $\mathbf{w}^\top \mathbf{x} + b = 0$, where $|\mathbf{w}^\top \mathbf{x}_n + b| = 1$.



The vector \mathbf{w} is \perp to the plane in the \mathcal{X} space: Take \mathbf{x}' and \mathbf{x}'' on the plane. $\mathbf{w}^{\top}\mathbf{x}'+b=0$ and $\mathbf{w}^{\top}\mathbf{x}''+b=0$,



The distance between \mathbf{x}_n and the plane $\mathbf{w}^\top \mathbf{x} + b = 0$, where $|\mathbf{w}^\top \mathbf{x}_n + b| = 1$.



The vector \mathbf{w} is \perp to the plane in the \mathcal{X} space: Take \mathbf{x}' and \mathbf{x}'' on the plane. $\mathbf{w}^{\top}\mathbf{x}'+b=0$ and $\mathbf{w}^{\top}\mathbf{x}''+b=0$, $\Longrightarrow \mathbf{w}^{\top}(\mathbf{x}'-\mathbf{x}'')=0$.

◆□▶ ◆□▶ ◆ ≧▶ ◆ ≧▶ ○ ≧ - の � で 5/20



The distance between \mathbf{x}_n and the plane $\mathbf{w}^\top \mathbf{x} + b = 0$, where $|\mathbf{w}^\top \mathbf{x}_n + b| = 1$.



The vector \mathbf{w} is \perp to the plane in the \mathcal{X} space: Take \mathbf{x}' and \mathbf{x}'' on the plane. $\mathbf{w}^{\top}\mathbf{x}'+b=0$ and $\mathbf{w}^{\top}\mathbf{x}''+b=0$, $\Longrightarrow \mathbf{w}^{\top}(\mathbf{x}'-\mathbf{x}'')=0$.

< □ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ = うへで 5/20



< □ ▶ < ፼ ▶ < 토 ▶ < 토 ▶ ○ ♀ ○ 6/20

... and the distance is...



Distance between \mathbf{x}_n and the plane: Take any point \mathbf{x} on the plane.



... and the distance is...



Distance between x_n and the plane: Take any point x on the plane.

Projection of $\mathbf{x}_n - \mathbf{x}$ on \mathbf{w} .

$$\widehat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \Rightarrow \mathsf{distance} = |\widehat{\mathbf{w}}^{\top}(\mathbf{x}_n - \mathbf{x})|.$$

4 ロ ト 4 日 ト 4 王 ト 4 王 ト 王 の 9 9 6/20



... and the distance is...



Distance between x_n and the plane: Take any point x on the plane.

Projection of $\mathbf{x}_n - \mathbf{x}$ on \mathbf{w} .

$$\widehat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \Rightarrow \mathsf{distance} = |\widehat{\mathbf{w}}^\top (\mathbf{x}_n - \mathbf{x})|.$$

distance =
$$\frac{1}{\|\mathbf{w}\|} |\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top \mathbf{x}| \Longrightarrow \frac{1}{\|\mathbf{w}\|} |\underbrace{\mathbf{w}^\top \mathbf{x}_n + b}_{=1.} - \underbrace{\mathbf{w}^\top \mathbf{x} - b}_{=0.}| = \frac{1}{\|\mathbf{w}\|}.$$

Point on the plain



Support Vector Machines - The optimization problem

Maximize the margin:

maximize_{**w**,b} $\frac{1}{\|\mathbf{w}\|}$ \implies Hard to solve subject to $\min_{n=1,2,\ldots,N} |\mathbf{w}^{\top} \mathbf{x}_n + b| = 1.$ We need to get rid of the min. Notice: $|\mathbf{w}^{\top}\mathbf{x}_n + b| = y_n(\mathbf{w}^{\top}\mathbf{x}_n + b).$ \mathbf{x}_n is classified correctly. minimize_{w,b} $\frac{1}{2}$ w^Tw \implies Equivalent problem subject to $y_n(\mathbf{w}^\top \mathbf{x}_n + b) > 1$ for n = 1, 2, ..., N;



Support Vector Machines - Constrained optimization

minimize_{**w**,b} $\frac{1}{2}$ **w**^T**w** subject to $y_n(\mathbf{w}^T\mathbf{x}_n + b) > 1$

subject to
$$y_n(\mathbf{w}^{\top}\mathbf{x}_n + b) \ge 1$$
 for $n = 1, 2, ..., N$,
 $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}.$

Lagrange? inequality instead of equality constraints \implies Karush-Kuhn-Tucker (KKT): Lagrange under inequality constraints



FSAN/ELEG815

Support Vector Machines - We saw this before Remember regularization?

minimize
$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^{\top} (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

subject to $\mathbf{w}^{\top}\mathbf{w} \leq C$.

Condition for the solution:

 ∇E_{in} relates to constraint. ∇E_{in} parallel to \mathbf{w}_{reg} but in the opposite direction.



 $E_{in} = \text{const.}$





Support Vector Machines - Lagrange formulation

minimize
$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} - \sum_{n=1}^{N} \alpha_n \underbrace{(y_n(\mathbf{w}^{\top} \mathbf{x}_n + b) - 1)}_{\substack{\text{constrain}\\y_n(\mathbf{w}^{\top} \mathbf{x}_n + b) - 1 \ge 0}$$

Why is it negative? Constraint is greater than or equal than zero.

w.r.t to w and b and maximize w.r.t each $\alpha_n > 0.$

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{n=1}^{N} \alpha_n y_n \mathbf{x}_n = \mathbf{0}$$

$$\frac{\partial \mathcal{L}}{\partial b} = -\sum_{n=1}^{N} \alpha_n y_n = 0$$

Note: $\alpha_n = 0$ $\frac{\partial f(\alpha_n)}{\partial \alpha_n} = 0$ Minimum of $f(\alpha_n)$ with $\alpha_n > 0$ Not in $\frac{\partial f(\alpha_n)}{\partial \alpha_n}$ but when $\alpha_n = 0$ ◆□▶ ◆□▶ ◆ ■▶ ◆ ■ ・ ○ へ ○ 10/20



◆□ ▶ ◆ ⑦ ▶ ◆ Ξ ▶ ◆ Ξ ▶ Ξ ジ ۹ ℃ 11/20

Support Vector Machines - Lagrange formulation Substituting

$$\mathbf{w} = \sum_{n=1}^{N} \alpha_n y_n \mathbf{x}_n$$
 and $\sum_{n=1}^{N} \alpha_n y_n = 0$

In the Lagrangian:

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



Support Vector Machines - Lagrange formulation Substituting

$$\mathbf{w} = \sum_{n=1}^{N} \alpha_n y_n \mathbf{x}_n$$
 and $\sum_{n=1}^{N} \alpha_n y_n = 0$

In the Lagrangian:

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

we get:

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

$$= \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} - \sum_{n=1}^{N} \alpha_n y_n \mathbf{w}^{\top} \mathbf{x}_n + \sum_{n=1}^{N} \alpha_n$$



Support Vector Machines - Lagrange formulation

Substituting

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

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



Support Vector Machines - Lagrange formulation

Substituting

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

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

we get:

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} y_n y_m \alpha_n \alpha_m \mathbf{x}_n^{\top} \mathbf{x}_m.$$

Now maximize w.r.t α subject to $\alpha_n \ge 0$ for n = 1, ..., N and $\sum_{n=1}^N \alpha_n y_n = 0$.

<□ ▶ < @ ▶ < E ▶ < E ▶ E の Q @ 12/20



Support Vector Machines - The solution

Notice: $\max \mathcal{L} = \min -\mathcal{L}$.

Quadratic programming:





Support Vector Machines - QP hands us α

Solution: $\alpha = \alpha_1, \alpha_2, ..., \alpha_N$ $\implies \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n.$

Many α_n equal to zero.

KKT condition: For n = 1, ..., N

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



Support Vector Machines - QP hands us α

Solution: $\alpha = \alpha_1, \alpha_2, ..., \alpha_N$ $\implies \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n.$

Many α_n equal to zero.

KKT condition: For n = 1, ..., N

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

α_n = 0 (interior points, no need for regularization)

 $E_{in} = \text{const.}$ \mathbf{w}_{lin}

We saw this before!

<□ ▶ < @ ▶ < E ▶ < E ▶ E の Q @ 14/20



Support Vector Machines - QP hands us α

Solution: $\alpha = \alpha_1, \alpha_2, ..., \alpha_N$ $\implies \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n.$

Many α_n equal to zero.

KKT condition: For n = 1, ..., N

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

- α_n = 0 (interior points, no need for regularization)
- $(y_n(\mathbf{w}^\top \mathbf{x}_n + b) 1) = 0$ (boundary points that support the plane)

$$\alpha_n > 0 \implies \mathbf{x}_n$$
 is a support vector.

 $E_{in} = \text{const.}$ \mathbf{w}_{lin} $\mathbf{w}^{\mathsf{T}}\mathbf{w} = C$

We saw this before!



Support Vector Machines - Support vectors

Closest \mathbf{x}_n 's to the plane. Support vectors \implies achieve the margin.





Support Vector Machines - Support vectors

Closest \mathbf{x}_n 's to the plane. Support vectors \implies achieve the margin.

$$\implies y_n(\mathbf{w}^\top \mathbf{x}_n + b)) = 1.$$

$$\mathbf{w} = \sum_{\mathbf{x}_n \text{ is SV}} \alpha_n y_n \mathbf{x}_n.$$

Solve b using any support vector:

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b)) = 1.$$





Support Vector Machines - Nonlinear transformation

 $\mathbf z$ instead of $\mathbf x$

$$\mathcal{L}(\alpha) = \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} y_n y_m \alpha_n \alpha_m \mathbf{x}_n^{\top} \mathbf{x}_m.$$

 $\mathcal{X} \longrightarrow \mathcal{Z}$







Support Vector Machines - Nonlinear transformation

 $\mathbf z$ instead of $\mathbf x$

$$\mathcal{L}(\alpha) = \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} y_n y_m \alpha_n \alpha_m \mathbf{z}_n^{\top} \mathbf{z}_m.$$

 $\mathcal{X} \longrightarrow \mathcal{Z}$







Support Vector Machines - "Support vectors" in $\mathcal X$ space

Support vectors live in the ${\mathcal Z}$ space. In the ${\mathcal X}$

space, "pre-images" of support vectors.





Support Vector Machines - "Support vectors" in $\mathcal X$ space

Support vectors live in the ${\mathcal Z}$ space. In the ${\mathcal X}$

space, "pre-images" of support vectors.

The margin is maintened in the \mathcal{Z} space.





Support Vector Machines - "Support vectors" in $\mathcal X$ space

Support vectors live in the ${\mathcal Z}$ space. In the ${\mathcal X}$

space, "pre-images" of support vectors.

The margin is maintened in the \mathcal{Z} space.

Generalization result

$$\mathbb{E}\left[\underline{E_{out}}
ight] \leq rac{\mathbb{E}\left[\# \text{ of SV's}
ight]}{N-1}$$





Example: Scikit-Learn on Synthetic Datasets

In high-dimensional spaces, data can more easily be separated linearly and the simplicity of classifiers such as linear SVM and SVM with a Radial Basis Function (RBF) kernel might lead to better generalization compared to other classifiers.



- Training points in solid colors.
- Testing points semi-transparent.
- Classification accuracy on the test set (lower right)



FSAN/ELEG815

Example: Scikit-Learn on Synthetic Datasets

