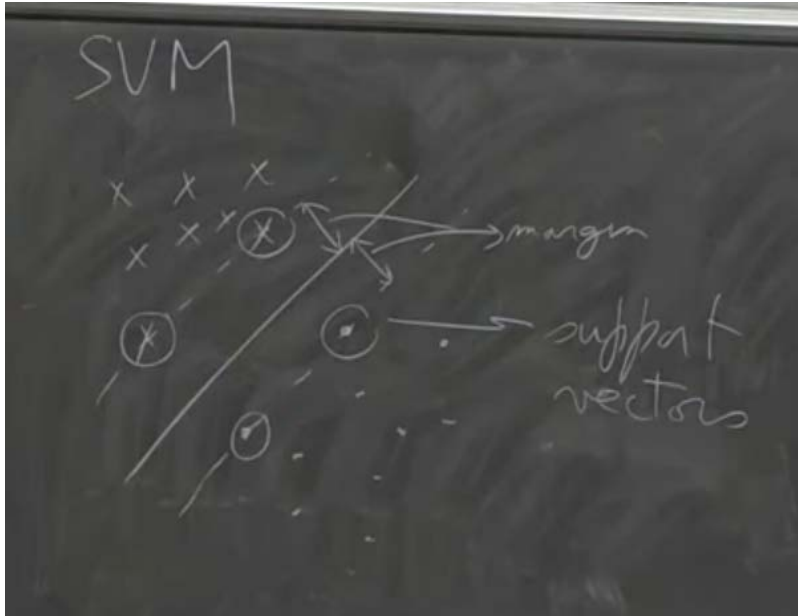


هنتكلم انهارده علي ال SVM .. الكيرنال ميثودز اللي بتشتغل علي الكلاسيكيشن .. زي ال SVM .. مهم جداً .. بس النويرال نتورك دخلت بعدو خلصت الدنيا خلاص .. لما اتكلمنا علي الكيرنال ميثودز بالنسبة للسكيلاتي هي cubic في حجم الداتا اللي عندك ... هل احنا لازم نعمل الميثودز ديه بتعتمد علي كل الداتا و لا هنقدر نشتغل علي جزء بس من الداتا .. ال SVM هو نوع من الكيرنال ميثودز ونبعتبره sparse و بيشتغل علي جزء من الداتا و بيرمي الباقي .. قبل ما نخش في الرياضه ونعقد الدنيا ... تعال نشوف نرسم شويه كذا ....

عندنا 2 كلاسيز من الداتا ... عاوزين نجيب linear separator يفصل ال 2 كلاسيز ... الطريقة الي هنفصل بيها هو اننا نقيس المسافه لأقرب نقطه والمسافه ديه هي المارجن .. اصغر مسافه لأقرب نقطه و بعدين نعمل نفس الحوار للناحية الثانيه .. ال SVMs بتلاقي linear sepearator line ف النص ... و النقط اللي بتبقا علي خط المارجن بتبقا هي support vector .. النقط هي عباره عن vectors و دول بيدولك linear separator لل decision boundary ..



ال SVMs هي برضو اسمها max margin classifier عشان هي بتلاقي linear line بيماكسيميز المسافه لأقرب داتا بوينت .. يعني الأوبتيمايزيشن بروبلم هتبقا هدفها انها ت maximize the margin ... ليه هنختار الماكس مارجن .. عشان الداتا ممكن تبقا نويزي ... فلو هرسم خط بيفصل .. وفي نقط فعلاً قريبه من الخط .. بس لو عندك إيزور فيالتالي النقط ديه ممكن تبقا في الناحية الثانيه من الخط .. بس لو عندك wide margin .. ففي فرصه احسن انك تخلي النقط اللي هتغلط فيها في الناحية الصح .. عشان نحسب ال distance بتاعت المارجن ... أول حاجه احنا عندنا لينير سيباراتور .. ده اللي هو  $wT$   $\phi(x)$  ... الدوت بروضكت بيطلعك يا بوزيتيف يا نيجاتيف علي حسب الدايركشن بتاع ال 2 فيكتورز ... فيالتالي لو عندك نقطه رايحه لكلاس واحد اديها 1 ولو رايح للكلاس الثاني اديها -1 .... فالفورميلا عشان تحسب المسافه بتاعت أي نقطه ل ال linear separator .. طيب اللي احنا عاوزينو هو اننا نلاقي شوية الويتس اللي هت maximize the margin .. المقصود بماكس المارجن هو المارجن اصلا هي اقل مسافه من النقط للخط .. فال minimum بالنسبة لل datapoint و دت الماكس بالنسبة للويتس عشان تلاقي الخط اللي يدبك أوسع مارجن ... هو ده الأوبجكتيف اللي هنعملو أوبتيمايزيشن

## Margin

- Linear separator:  $w^T \phi(x) = 0$

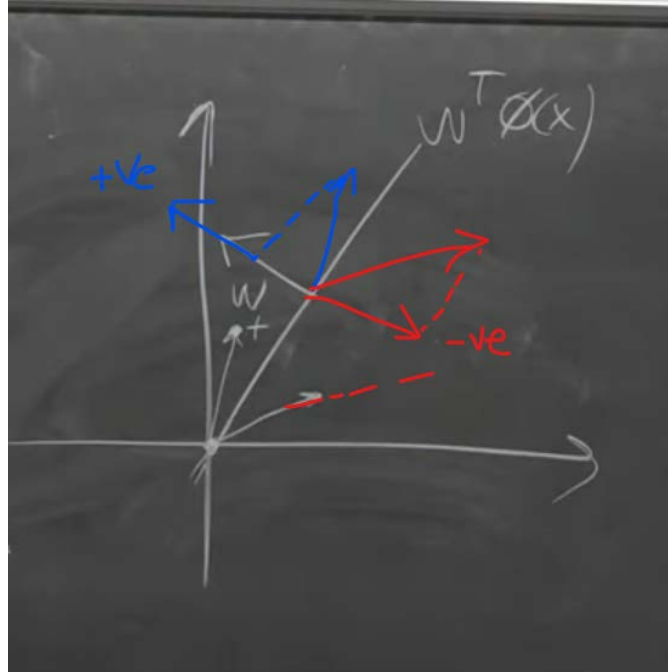
- Distance to linear separator:

$$\frac{y w^T \phi(x)}{\|w\|} \text{ where } y \in \{-1, 1\}$$

- Maximum margin:

$$\max_w \frac{1}{\|w\|} \left\{ \min_n y_n w^T \phi(x_n) \right\}$$

الدكتور هيتكم شويه علي البرسيبترون .... هو برضو بيلاقي لينير سيبيريتور ... بس ده بيعتمد علي ال starting values بتاعت الويتس ... فاللي بتوصلو في الآخر هيبقا مختلف علي حسب ال starting values .. المهم ان ال update rule هو انك هتجمع او هتطرح كل misclassified datapoint



الدكتور بيقول ان البرسيبترون prone to overfitting .. لو عندنا SVM احنا عندنا unique max-margin و عندنا linear separator .... ال update rule هي مش سهله ... و بتبقا عندك convex quadratic optimization و الدكتور بيقول انو roudust to overfitting ...

SVM	Perceptron
• unique max-margin linear separator	• linear separator (depend on starting values)
• convex quadratic optimization	• simple update rule
• robust to overfitting	• prone to overfitting

## Maximum Margin

- Unique max margin linear separator

$$\max_w \frac{1}{\|w\|} \left\{ \min_n y_n w^T \phi(x_n) \right\}$$

- Alternatively, we can fix the minimal distance to 1 and minimize  $\|w\|$

$$\min_w \frac{1}{2} \|w\|^2$$

s.t.  $y_n w^T \phi(x_n) \geq 1 \quad \forall n$

at least dist 1

- This is a convex quadratic optimization problem that can easily be solved by many optimization packages

طبيب دلوقت عاوزين نشتغل بال dual representation ... احنا مش لازم ندفع تمن لل new mapping phi .. انما هنستخدم new mapping phi .. وهنعمل ريمابنج للداتا من الأوريجنال سبيس لل new space .. وهناك كل الحسابات هتبقا in terms of dot product ... وتقدر تبدلها بكيرنال فانكشن .. ده اللي عملناه ل Gaussian processes و ديه كانت مجرد ال dual version بتاعت ال Bayesian linear regression وهنا مع ال SVM في برضو dual version ... في advantage للدوال ده .. ان ال optimization is sparse بشأن الداتا بوينتس اللي علي المارجن بس هي اللي تهمني .. الأوبتيمائيزيشن ثاني هي اننا نقلل ال w مع الوضع في الاعتبار الشرط اللي اشتغلنا عليه ...

المهم دلوقت اللي نقدر نعملو هو اننا نعيد كتابة الحوار ده بطريقة اننا ن move the constraint into the objective .... في ال optimization مش بيبقا convenient ان بيبقا عندنا constraint .. مثلا لو بصينا علي جريدنيت ديسنت .. وانت بتحسب direction بتاع ال steepest descent خطوه هناك .. لو عندك كونسترينت .. انت بتلصيح ال gradient و بتلصيح الأوبتيمائيزيشن فالدنا بتبقا معقده أكثر .. فشيل ال constraints اللي عندك بشأن الدنياي تبقا اسهل ... المهم انك ممكن تطالع ب equivalent optimization problem وهنا احنا هنشيل ال constraint ... هنحطه في ال objective ك penalty term ..... وهنا بنعمل lagrangian وديه new objective اللي هي هتبقا ماسكه الترم بتاع المنيم عادي .. انما ال constraints دخل علي هيئة penalty .. هنتقيس المسافه اللي عندك .. لو مش أكبر من واحد فهتخط penalty و الماجنتيود بتاعها هيبقا an ... لو المسافه كان أكبر من الواحد .. حط ال an = 0 بشأن مييقاش عندنا أي penalty ...

## Dual derivation

- Transform constrained optimization

$$\min_w \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_n w^T \phi(x_n) \geq 1 \quad \forall n$$

into an unconstrained optimization problem

- Lagrangian

$$\max_{a \geq 0} \min_w L(w, a)$$

$$\text{where } L(w, a) = \frac{1}{2} \|w\|^2 - \sum_n a_n [y_n w^T \phi(x_n) - 1]$$

penalty for violating the  $n^{\text{th}}$  constraint

newname  
10/14/2021 8:16:46 AM

هنحل كلورد فورم ال inner minimization ... ده بيحصل في كلورد فورم .. ديه convex quadratic objective ... خذ بالك ان quadratic هنا عشيا ال L2Norm .. الكونفكس معناه ان عندك جاوبال اوبتيمال منيمم ... المهم في الآر وصلنا ال w هتبقا عامله ازاي ... بقت عبارة عن لينير كومباينيشن من الداتا بوينتس في النيو فينشر سبيس .. ال an هنا مش بتساوي صفر .. والداتا بوينتس هنا هم السابورت فيكتورز .. فال w هتبقا linear combination of support vector



- Solve inner minimization:  $\min_w L(w, a)$

$$\min_w \frac{1}{2} \|w\|^2 - \sum_n a_n [y_n w^T \phi(x_n) - 1]$$

- Set derivative to 0

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_n a_n y_n \phi(x_n)$$

- Substitute  $w$  by  $\sum_n a_n y_n \phi(x_n)$  in  $L(w, a)$  to obtain:

$$L(a) = \sum_n a_n - \frac{1}{2} \sum_n \sum_{n'} a_n a_{n'} y_n y_{n'} k(x_n, x_{n'})$$

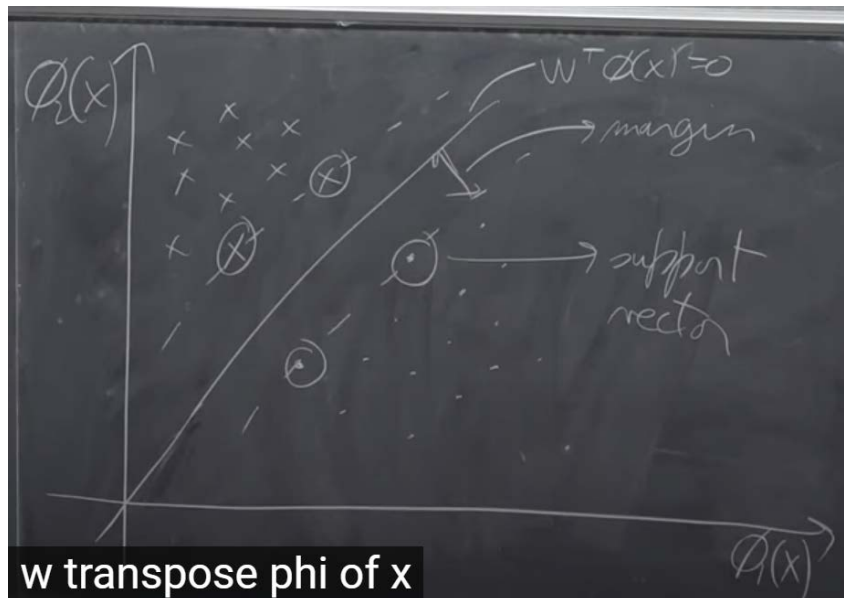
# Dual Problem

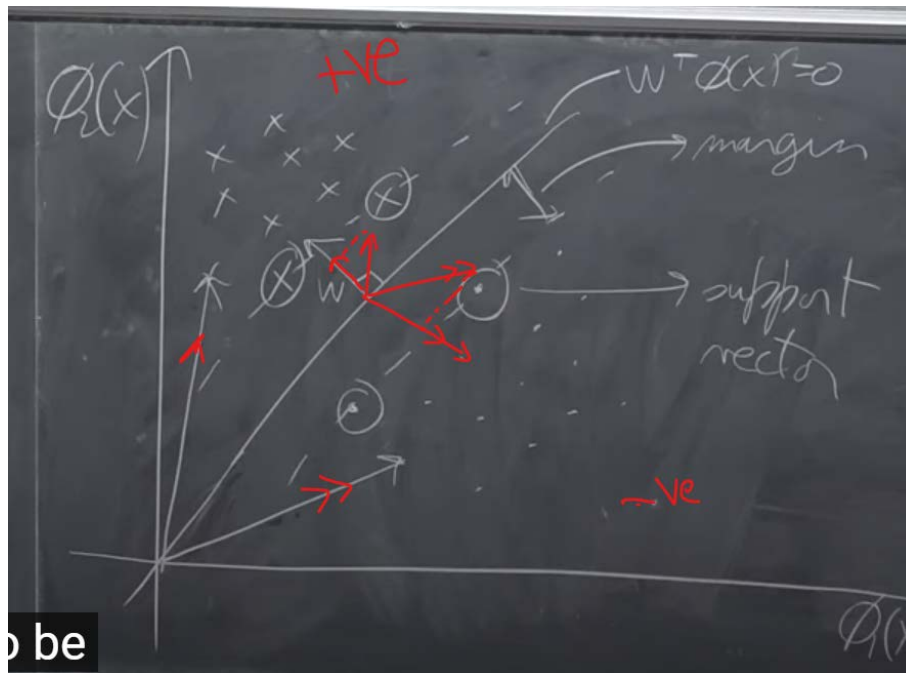
- We are then left with an optimization in  $\mathbf{a}$  only known as the **dual problem**

$$\begin{aligned} \max_{\mathbf{a}} L(\mathbf{a}) \\ \text{s.t. } a_n \geq 0 \end{aligned}$$

عنا  
نقوم  
بالتحسين  
في  $\mathbf{a}$

- **Sparse optimization:** many  $a_n$ 's are 0





## Classification

- Primal problem

$$y_* = \text{sign}(w^T \phi(x_*))$$

*features*

- Dual problem

$$y_* = \text{sign}\left(\sum_n a_n y_n \phi(x_n)^T \phi(x_*)\right)$$

$$y_* = \text{sign}\left(\sum_n a_n y_n k(x_n, x_*)\right)$$

*sum of degree for similarity*

newname  
10/14/2021 9:25:13 AM

kernel functions captures the similarity between pairs of data point ... it's dot product into feature space.. If data points are similar, their vectors are going to point more or less in the same direction, there for the kernel function will be positive. If not similar, the kernel function is negative ...