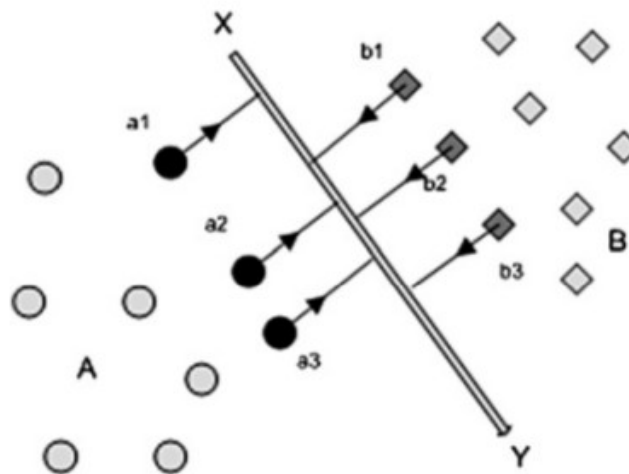


From online tutorial paper:

Support vector machines (SVMs) SVMs, a discriminative learning approach, classify inputs (eg, words) into categories (eg, parts of speech) based on a feature set. The input may be transformed mathematically using a 'kernel function' to allow linear separation of the data points from different categories. That is, in the simplest two-feature case, a straight line would separate them in an X-Y plot: in the general N-feature case, the separator will be an (N-1) hyper-plane. The commonest kernel function used is a Gaussian (the basis of the 'normal distribution' in statistics). [The separation process selects a subset of the training data \(the 'support vectors' - data points closest to the hyperplane\) that best differentiates the categories.](#) The separating hyperplane maximizes the distance to support vectors from each category.



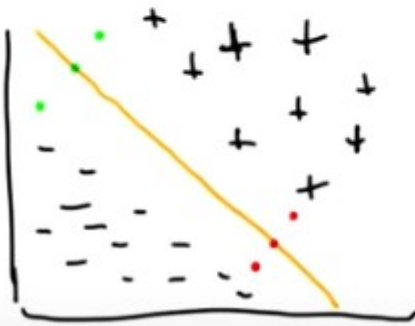
正常 note:

1. Hi Michael, how are you doing? >> Oh, good thanks. What the, what's on tap for today? >> Well, as you can see on your screen today we're going to talk about support vector machines. >> Hmm. >> And at the end of it, we are going to circle back to boosting, so that I can try to explain to you why it doesn't seem to over fit as much as you might think or at least not in the ways that you make it. >> Ooh, that would be cool. I was worried that that was maybe a command. >> What was a command? >> The thing that you wrote on the screen. >> [INAUDIBLE] [LAUGH] Man you know that reminds me of one of my favorite little thing about police. That any combination of the word police is a legal sentence >> [LAUGH] I see. >> Please. >> Please. Police Please please. Please police. >> Please please please please please please police. Those are all legal sentences. >> All right, we'll please stop. >> [LAUGH] Man, this is why I hang out with you, Michael. Ok,so today we're going to talk about support vector machines and I'm going to do something unexpected. I'm going to start out The beginning of this with a quiz. >> With a quiz. >> Yes, with a quiz. So here's the quiz for you, Michael. I've drawn on here ,uh, some points, some labeled positive and some labeled minus, representing two different classes. And, you'll notice that you could draw a line to separate them, therefore they are >> Linearly separable. >> Exactly. So the question I have for you is simply this, which is the best line? And here's how you're going to show me. Here's how you're going to show me. I'm not just going to ask you to draw some random line because that's too hard for us to get feedback on. So instead what I've

done is I've drawn three green dots. Here on the sort of upper left and three red dots here on the bottom right and what I want you to do is select one of the green dots and one of the red dots and since you'll have two dots that will define a line and that'll be the way you indicate to me which of the lines is in fact the best line. >> The best line. >> The best line. And I'm not going to even tell you what best means. You tell me what best means, when you justify why you would choose one over the other. >> If I think they're all the best, I can choose any one I want? >> That's right. Ok. So you got it? >> I think so. >> OK. Go.

SUPPORT VECTOR MACHINES

QUIZ!

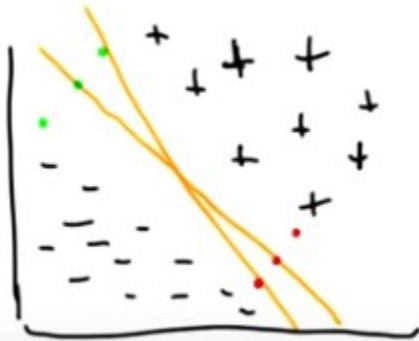


WHICH IS THE BEST LINE?

2. Alright, Michael. So, you got it? Which line do you think is best? >> Alright, so if I choose one green and one red [CROSSTALK] that means there's nine different possible lines. But they all separate the points. >> That is correct. >> I'm pretty sure. I, you know, I think they're intended to even if they don't quite. But I think they do, I think they actually all separate the positives from the negatives [CROSSTALK]. The positives on the upper right side and the negatives on the lower left side. So, I mean the only one that seems further special is the middle one. If I choose the middle green and the middle red it's kind of, you know, ecstatically pleasing, there's a lot of space on, on each side. >> Okay. >> So, I'm going to go with that but it's not clear, you know again, best if I was a photographer that might be the best. >> Hm, you are a photographer. So, let's pretend I can draw a straight lines. That's pretty good. That's the line that you've chosen. >> Yes. >> And you think that's best. Well, I'm going to give you a hint, Michael, and tell you that you were correct. That's a pretty good hint. >> Yeah.

SUPPORT VECTOR MACHINES

QUIZ!

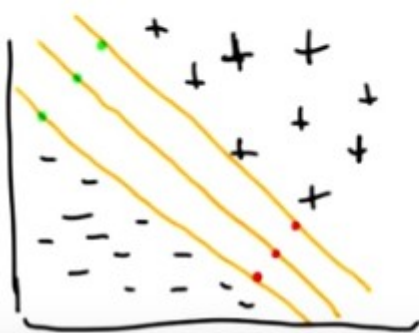


WHICH IS THE BEST
LINE?

>> But now I want you to figure out for me why that line is better than this line. >> Interesting. Alright, well so one thing that's, that second line. The orange line(即剛加的那條), as oppose to the orange line, has against it, in some sense it seems to get really, really close to that bottom minus point. Maybe it's a little too close. Like maybe that minus is near other minuses that we just can't see. And since we drew really really close, a line really really close to it, it could be we ended chopping of some of the minuses that. That we might want to have kept on the other side. And so maybe by drawing it in the middle it's sort of, we have the biggest, I don't know like demilitarized zone. >> No, I like that. So that is a very good explanation for why you would prefer the middle line over the. Other line that isn't the middle line. And you can make a similar argument for any other of the nine lines that, that or the other eight lines that you could possibly draw.

SUPPORT VECTOR MACHINES

QUIZ!

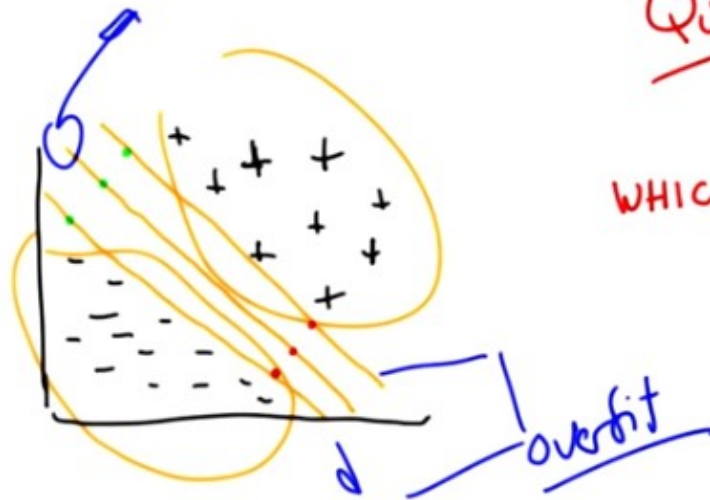


WHICH IS THE BEST LINE?

Let me draw another line for you and you tell me if you think this one is good or bad and why it might be better or worse than the middle one. See if you give me a different answer. I can take one of the parallel lines which again we will pretend is a straight line and put it there. Or I can take in fact the other line that is meant to be parallel to it and why aren't those lines better? Or just as good. >> I guess I'd use the same explanation as before which is that these, the, the line that was, is really close to the pluses maybe is a little too close for comfort and it, it gets very close to that plus. The, kind of, I don't know. I could point to it, but maybe you should point to it. The plus that's really close to the line. >> I'm pointing to it. >> Because again if there's a just, you know, other data points near there, it's going to make this distinction between the positives and the negatives that isn't really warranted by the data. >> Right, so I could make the opposite argument with you, Michael, which is that listen, they all separate the points. They completely explain the data as we see them. So, all three of those lines explain the data. In fact, all nine of those lines explain the data. So, why aren't they warranted by the data. What's the problem that you might run into if you put one line very close to the positives or one line very close to the minuses? >> Well, again...so, so...it may be that they fit this data. But, this is just a sample from the population. And so... You know, we don't know what's going to happen really close to that other plus is. So like in the nearest neighbor algorithm for example which you always try to make me remember, it's going to want to, you know, if it's closer to the pluses, it should be a plus.

SUPPORT VECTOR MACHINES

QUIZ!

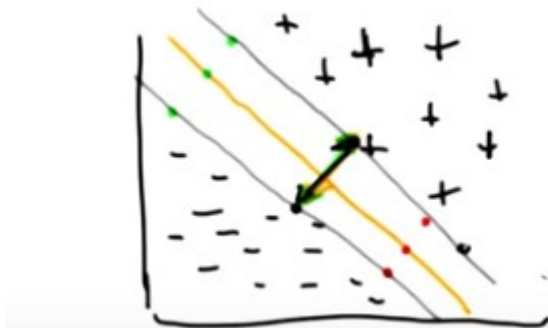


WHICH IS THE BEST LINE?

>> Right. That makes sense and I think that's exactly right so your intuition and I think the intuition that people should take with them is that lines that are too close to the positives, or lines that are too close to the negatives, sort of have this feature where you're believing 'the training data that you've got' too much. You've decided that all of these boundaries are fine because of the specific set of training data that we've got. And while you want to believe the data, because that's all you've got, you don't want to believe the data too much. >> That's overfitting. >> That's overfitting. So, the problem with those two lines, or one way to think about the problem with those two lines, is that these lines are more likely to overfit. And why is that? It's because they're believing the data too much. So, given that we're trying to avoid overfitting, what could you say about the middle line? So, here's the argument that I would make, and you tell me if you buy it, Michael. This line, or the line it's intended to represent anyway, is the line that is consistent with the data while committing least to it. That seems like a good way of saying what I was feeling, yeah. It is funny though because like, overfitting up to this point, we you, we were generally talking about overfitting as being something where there's a great deal of model complexity in some sense. >> Mm-hm. >> And it doesn't seem like those lines that are closer to the pluses or closer to the minuses. Are inherently more complex, they're still lines. It's interesting that they, they, they kind of maybe behave as if they are. >> Right, and in fact they, they're, it's a more, sort of, literal interpretation of the words over and fit, right? You, you have decided to fit the data, and you believe it too much, and what you really want to do is commit. The least that you can commit to the data while still being consistent with it right. So this basic idea of, uh,uh, finding the line of least commitment in the linear separable set of data, is the basis behind support vector machines. So what I want to do next is I want to see if we can come up with some equation that would help us define such a line. It's easy in this case cause we're staring at it but if you imagine these points were in 700 and. Thirteen thousand dimension it would pretty difficult to find such a plane just by staring at it so let's see if we can try to work out how you go about finding this least commitment line okay >> Cool.

SUPPORT VECTOR MACHINES

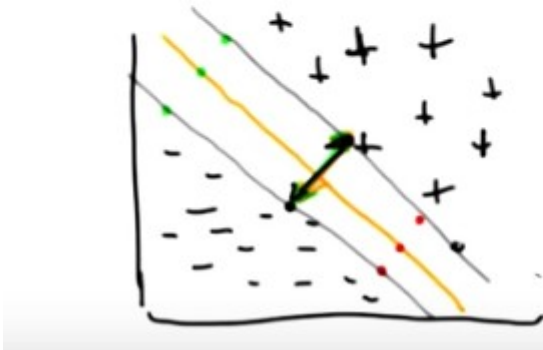
$$y = mx + b$$



3. Okay Michael. So, let's, [let's write down a few equations](#). Let's try to be a little bit more formal. A little bit more mathematical about this idea. So, I'm going to try to encapsulate, what we just talked about, to the line of least commitment. By drawing another line. So, if we think, of this top gray line, here. As sort of the line, that gets as close to the plus points as possible, without crossing over, to them, and mis-classifying them. And we think of the bottom gray line, as the one that gets as close as possible, to the minus signs, without crossing over them and, and giving mis-classification. And then the middle line, is sort of in the happy medium. Okay? So, [what you really want from the lines](#), all the lines that are possible, between these two boundary lines. If I can use that language. [Is that](#), somehow, [this distance](#)(即兩條灰線間之距離) here, [is as big as possible](#). Can you see that? >> Yeah, though it seems like the gray line, could be pushed out a little more, right? because the minuses don't bump into it. That's a good point Michael and I'm going to fix that, by putting a minus sign here. >> [LAUGH] >> Okay [CROSSTALK] >> I see, this is [CROSSTALK] >> I did the best I could under the circumstances. >> Data, revisionist history. >> No, there's just an invisible point, and I just made it more visible. For the sake of the reader. Okay, so, I've got these two lines, these are sort of as far as I can go, without stating to do mis-classification with my, my separating line, and the line in the middle, we've already argued is the sort of the best one because it provides the least commitment. So, that means, you want to have a line, such that, it leaves as much space as possible, from the boundaries. Alright, Michael. So let's see if we can figure out exactly, what that line is like. So, the first thing, that I want to do is, is introduce a little bit of notation. Right? [So we all remember what the equation, of a line is](#). It's $y = m x + b$, that's just a general equation for a line. But, here even though we're going to be drawing with lines, we really want to deal with the general case, where we're talking about hyperplanes.

SUPPORT VECTOR MACHINES

$$y = mx + b$$
$$\underline{y} = \underline{w}^T \underline{x} + \underline{b}$$



$y = mx + b$ 和 $y = w^T x + b$ 之區別:

$y = mx + b$ 為 x - y 坐標係中的一條直線, 在這條直線之上的點為+, 在這條直線之下的點為-.

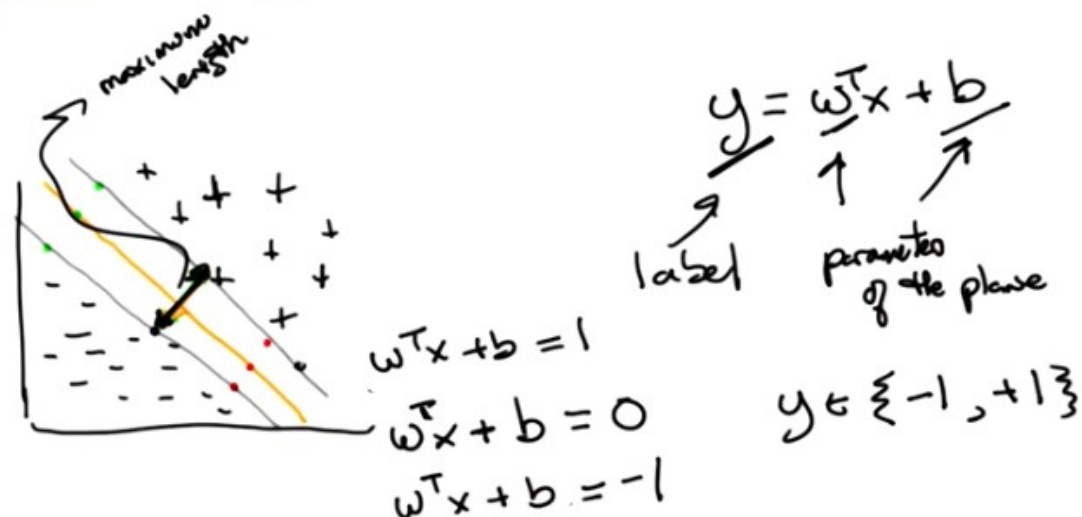
$y = w^T x + b$ 中的 x 是一個向量, 它即多維空間中的一個點的坐標(故 x 為 input). y 是 output, 後面說了, 它的值只能取 -1 和 +1 兩個. 對於某個點 x , 用 $w^T x + b$ 算出來的值 若為 +1 (即 $y = +1$), 說明此點為+, 若為 -1 (即 $y = -1$), 說明此點為-.

後來由後面知, 更準確的是: 若用 $w^T x + b$ 算出來的值 大於 +1, 則說明此點為+, 若 小於 -1, 則說明此點為-.

由後面馬上即可知: $y = mx + b$ 即 $w^T x + b = 0$ 的二維形式.

And generally, when we write about hyperplanes, we describe them as some output, let's just call it $y = w^T x + b$. here, because of what we're, what we're trying to do with classification, the output y is going to be some value that indicates, whether, you're in the positive class or you're in the negative class. w are the parameters, or w represents the parameters for our plane. Along with b , which is what moves it out of the origin. Okay, are you with me. >> I think so, but that's, so may be we should get rid of that top Y , because, that Y is different kind of Y . Alright so the top Y is talking about the Y dimension of the plane and in the second equation, that Y is kind of folded into the X . And we have a new y , which is actually, the output of the classifier. >> Right. I like it, so let's get rid of that first y which is just an equation for a line and let's ask what each of these things are.

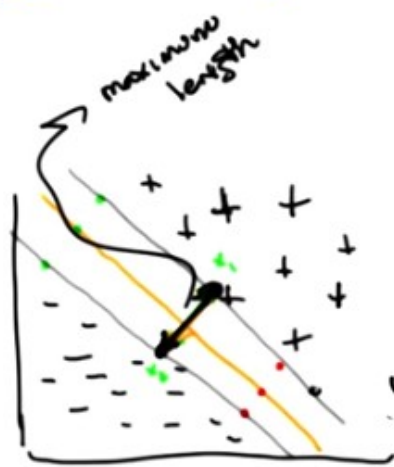
SUPPORT VECTOR MACHINES



So, let's just say that again for clarity's sake. I think, you make a good point, Michael. **y here, is going to be our classification label**, right, whenever we're talking about using a linear separator, effectively, what we've been talking about, which I realize now we never ever actually said explicitly, is that you are taking some new point, projecting it onto the line, and then looking at the value that comes out from projecting it. And in this case, in particular, **we want positive values to mean yes, you are part of the class, and negative values to mean that you aren't a part of the class.** Okay? >> Yeap. >> This is our classification label y. W represents again, the parameters of the plane. Along with b, which is what moves it in and out of the origin. So, this is now, effectively, what our linear classifiers actually look like. Even in multiple dimensions, with hyperpoints. Okay? Cool. So, let's take that, and and, and push it, to, to sort of the next level. Let, let's figure out exactly, what we would expect the output, of our hyperplane, to be in this example, that I, I've drawn on the screen here. So, we know we want to find this orange line in the middle, which has the property, that, it is your decision value, it tells you whether you are in the positive class or negative class, on the one hand, but also, that it has the property of being, as far away, from the data as possible, while still being consistent with it. **So, if you're on the decision boundary, for this particular line, which again, is $w^T x + b$.** What would be, the output, of this classifier for any point that lies along the line? >> so, right. **If that's decision boundary, that's where it's kind of not sure if it's positive or negative, so that should be zero.** >> Right, so the equation of this line or this hyperplane, is $w^T x + b = 0$ (由於 x 為一個點的坐標, 故 $w^T x + b = 0$ 表示多維空間中的一個平面, 或二維空間中的一條線, 上圖中的例子為二維空間, 故 $w^T x + b = 0$ 表示的是中間那條黃線. 前面的 $y = mx + b$ 即 $w^T x + b = 0$ 的二維形式). For, some set of parameters w and b, we don't yet know what they are. But, what we do know, that this is the definition of a hyperplane. Where the equation for a hyperplane, and since it's at the decision boundary, it should give me neither a positive or a negative output. Okay? >> Yep. >> Okay now, one question we can ask ourselves then, **if we look at these**

other lines as well, what's the equation for the other grid lines? That are right at our positive or negative examples. So to help you answer, that, I want to talk about what the labels themselves ought to be. So, just like we did with boosting, let's say that our labels are always going to be, from the set minus one and plus one. We know that our labels, are -1 and +1. And so we're going to take advantage of that fact, by saying well. The line that brushes up, against the positive example. We want it to be the case, since, we know those labels are going to be plus 1. That, the output, of that particular line, that particular linear separator, would be plus 1, on the very first point, that it encounters. Does that make sense? >> Yeah >> Okay >> That way the things, that are, that it's you know that are kind of past the line are going to be above 1 and the things are before the line, in that kind of demilitarized zone, are going to be between zero and 1. >> Right, so in fact given what you just said, what is the equation of that line? >> Oh I see. So, it should, be the $w^T x + b = 1$ for the top gray line (為何是 $w^T x + b = 1$? 因為我已總結出, 若用 $w^T x + b$ 算出來的值 大於+1, 則說明此點為+. 故 $w^T x + b = 1$ 就是+區之邊界). > > That's exactly right. And by a similar argument, where would you say the, the line of the hyperplane should be for the bottom gray line? Analogously, it seems like that one should be -1. >> Right. So, we, we have the decision boundary, we're looking at, and we know that the equation of the line is $w^T x + b = 0$. We know that if we slid that line, towards the positive values, we would end up with $w^T x + b = 1$. And if we slid it towards the negative values, we'd end up with, $w^T x + b = -1$, now we can ask ourselves, how this helps us. And it helps us in a very a simple way. We know that we want the boundary condition line, the one that is actually our decision boundary, to be as far as possible from both our positive and negative examples, so that would mean then, and I hope you buy this Michael, that the distance, between the two gray lines, which are parallel, to that line, needs to also be maximum. >> Yeah, that's exactly what we want. >> Right, so we want this vector (兩條灰線間的那個 vector) here, to have, the maximum length that we can have. Okay, so, how are we going to figure out how long, that particular line is. So, here is a simple idea. Well, the lines are parallel to one another. We can pick points on that line, to define, that particular distance there. So, just because, it's really easy to do the math, I'm going to chose a point here and a point here. Those points have the property that if I draw the line between them, you get a line that is perpendicular to the two gray lines. Okay? And I don't know what those x values are, but I do, I'm just going to call, them x_1 and x_2 (x_1 和 x_2 即那個 vector 的兩個端點, 即它們為常數, 跟用到 x_1 和 x_2 不同). That seem fair? >> Hm, I guess that seems as good a name as any. >> And that is going to define the two points that I have, and the distance between them, is in fact going to be or the, the vector, that is defined by their difference is in fact going to have the length that tells you how far apart those two (gray) lines are. Which in turn because of the way that we've constructed them tells you how far apart your boundary decision line (指中間那條黃線) is from the data and we want that to be maximal, because then we made the least commitment to the data (這就是為何要使兩條灰線間距盡量大 ← 因為我們想中間那條黃線 make the least commitment to the data).

SUPPORT VECTOR MACHINES



$$w^T x_1 + b = 1$$

$$w^T x_2 + b = -1$$

$$w^T x + b = 1$$

$$w^T x + b = 0$$

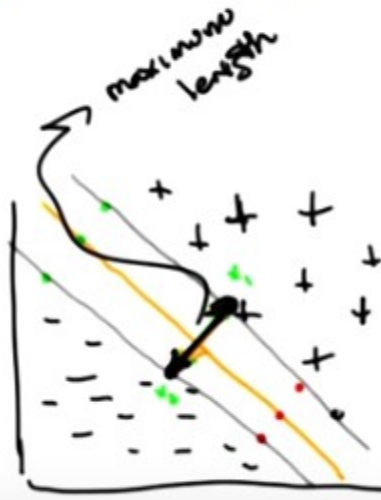
$$w^T x + b = -1$$

So, let's write that down as algebra. The equation for our positive line, so to speak, is therefore $w^T x_1 + b = 1$ (x_1 是上面那條灰線上的任意一點). And all I've done there is substitute, some point, I don't have to know what is, it's going to turn out, that gives me some point on that line, okay? It gives me, puts me in some particular place, on that line. And similarly, I can do the same thing, for my negative line. And get $w^T x_2 + b = -1$ (x_2 是下面那條灰線上的任意一點). Now, we want the distance between, these two, hyperplanes or these two lines, in this example, to be maximal. And in order to figure out what, that means, we need to know exactly, what that line is. So, it's, it's the difference between the two. So, we can just use our, our favorite trick, when we're doing systems of linear equations. And, just subtract the two lines. We basically have, two equations and two unknowns. And, we simply subtract them, from one another. So, that we can get a single equation. That happens to represent, the distance, between them. So, if I subtract the two from one another, what do I get? >> Quiz. >> Oh, I like that, we get a quiz. That is the correct answer. Even though it seems like a tightness match, in fact quiz, is the correct answer [LAUGH]

4. Okay. So here's the quiz. Michael is going to answer it, but we want to give you a chance to answer it first. I've got these two equations of two different hydroplanes, though they're parallel to one another, 'because they have the same parameters. That is to say, I have two equations and two unknowns. I want to subtract them from one another and what we want you to do is we want you to solve for the line that is described by their difference. Do you understand that, Michael? Or, have I, have I told them enough, or not? >> Yeah, I think I'm just going to subtract the second equation from the first equation. It seems

pretty straightforward. >> Okay, it seems reasonable to me. But remember, the output that I want you to figure out here is exactly what the distances between those two planes, okay? That is, between what's represented by X_1 and X_2 , okay? Go.

SUPPORT VECTOR MACHINES



Quiz

$$w^T x_1 + b = 1$$

$$w^T x_2 + b = -1$$



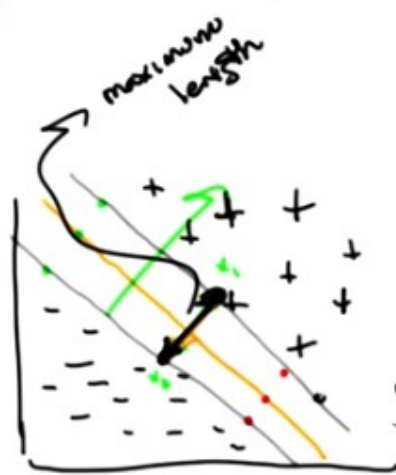
$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

$$w^T (x_1 - x_2) = 2$$

SUPPORT VECTOR MACHINES



$$\begin{aligned} w^T x + b &= 1 \\ w^T x + b &= 0 \\ w^T x + b &= -1 \end{aligned}$$

Quiz

$$\begin{aligned} w^T x_1 + b &= 1 \\ w^T x_2 + b &= -1 \end{aligned}$$

$$\frac{2}{\|w\|}$$

$$\frac{w^T (x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$$

margin

5. Okay Michael, what's the answer? >> Well, there's the answer to the question that I thought you were asking and then there's the question that you then at the end actually asked. It seems like at the end you asked what is the difference between the two, the distance between the two lines and I I feel that, like that's just The, like the norm of x_1 minus x_2 . But, the difference between these equations is going to be well ,uh, W transpose. >> Well why not write down what you're, what you're telling over in the side over here and then we can put the final answer in the box. >> Okay. >> Okay, so what now? >> W transpose. x_1 minus x_2 , equals two. >> Right, so you used ,ah, the power of subtraction to make that work. Okay, very good. Okay, so that's the difference between those two equations, now how am I going to go from there to figuring out the distance between x_1 and x_2 ? >> I still feel like its just that norm of x_1 minus x_2 . >> Okay, but I want you to tell it to me in terms of w . And what I want you to tell it to me in terms of w , because the only thing we have to play with here is w , well w and b . That's what defines our line. And so I want to find the right line, the only thing I get to play with, W and b . So, I'd like to know something about the relationship between w and the distance between x_1 and x_2 . >> Well, w transpose times their differences too. >> [LAUGH] That's true. >> That's not telling us the distance, though. So what is the distance in terms of W ? >> Well what if I told you w was a number? What would you do if it was just a simple scalar and you had this equation, and I wanted to know what x_1 minus x_2 was. What would you do? >> And I wanted it in terms of W ? >> Yeah, I wanted to know what x_1 minus x_2 was equal to. >> Oh, I see. So, if I divide it by w , that would be helpful' cus then x_1 minus x_2 would be two over w . >> Right, but you can't do that because w is a vector, and you can't really divide by a vector, at least not in the world that we're talking about. So how are you going to... Make that work. Would you like a hint? >> Sure. >> Well here's a hint ,we want to move w over from one side to the other. We could start doing all kinds of tricks with inverses, and with the inverse of a vector. There's all kinds of things that you could do, but actually the easiest thing way of doing it is getting rid of w on one side. And the easiest way to do that. Is to divide both sides by the length of W . So, rather than dividing both sides by W . We divide them by the length of W . Now what is dividing W by the length of W , give you? >> So, right. So W divided by the length of W , is a normalized version of W (即 W 方向的單位向量). So it's like Something that points in the same

direction as W , but sits on the unit's sphere. >> Right. No, that's exactly right! Alright, in fact it is a sphere because it's a, it's a hyper sphere I suppose. right! So we do that and that effectively is like giving you a value one, like you said it's a unit sphere. And so now we're actually talking about the difference between the vector X_1 and the vector X_2 projected onto the unit sphere and that's equal to two over the [CROSSTALK]. >> Wait. That's what $(x_1 - x_2)$ is projected onto the W vector, the normalized W vector. >> Right, the doubled norm, the normalized >> The double normalized U vector, right >> [LAUGH] So does that help you? >> Does that actually answer the question? Doesn't seem like it does. >> no, it does. >> So is that, okay, alright, hang on [LAUGH], so, the x_1 minus x_2 , dotted with W . So W , W , we don't know. It could be anything. Can it? So. >> Mm-hm. >> We've taken x_1 minus x_2 , projected it onto W . So it's like the, the length of x_1 minus x_2 , but in the w direction. >> Exactly. That's exactly right. So, the link, which, by the way, is, well, is exactly right. So, what we've just done is we have found the link of x_1 and x_2 in the W direction. What do we know about w with relationship to the line? >> W is the parameter to the line. >> Yes, but in particular. As with, in the case with any hyper plane, W actually represents a vector that's perpendicular to the line. >> And since we chose X_1 and X_2 , so that they would in fact, their distance or the difference, would in fact be perpendicular to the line. What we've just done is projected. That rose the difference between those two vectors onto something that is also perpendicular to the line. And so what that ends up giving us is, in fact, its length. So we maximize the length. Of x_1 minus x_2 with doing what with W . >> I see so the thing on the left is, in fact, have we answered the quiz yet by the way, or are we still working on that? >> I'm going to say we are still working on it. >> Oh men, alright this is a hard quiz. The thing on left, not just were the braces are, that actually turns out to be the distance between the two. Hyperplanes. >> Right, let's let's give that a letter. Let's call it M . >> Mm. >> Mm. >> And, we're saying that equals two over the norm of W . And that's, so, if we want to maximize that, the only thing that we have to play with is W and that is made larger and larger as W gets smaller and smaller, in other words, pushing it toward the origin. >> Right. >> So it set W s to all zeroes, and we should be golden. >> Right, except if we push all the W s to zero, we might not be able to correctly classify our points but what this does tell us is that we have a way of thinking about The distance of this vector and where the decision boundary ought to be. We want to find the parameters of the hyperplane (圖中的 那兩條灰線 和 那條黃線 即多維空間中的 hyperplane) such that we, maximize this distance over here represented by this equation while still being consistent with the data. Which makes sense because that's actually what we said in the first place. By the way, this thing has a name and it's. The reason why I chose M , it's called the margin, and what all of this exercise tells you is that your goal is to find a decision boundary that maximizes the margin, subject to the constraint that you actually want to correctly classify everything, and that is represented by that term. Now somehow, it feels like having gone through all this we outta be able to use it for something and turn in into some other problem we might be able to solve. Might be able to maximize for, uh, so that we can actually find the best line. And it turns out we can do that. >> Have we answered the quiz yet? >> Oh yeah we did. Which is in fact what I wanted so the answer is. >> Wow. Somebody gets that, that would be pretty impressive. >> That would be very impressive. Or anything similar to this I would accept. [LAUGH] In fact I probably better will. Okay good. So, it turns out that this whole notion of thinking about finding the optimal decision boundary is the same as finding a line (a line 是因為 那兩條灰線 和 那條黃線 的參數是一樣的, 故只找出 那條黃線 就可以了) that maximizes the margin. And we can take what we've just learned, where we've decided the goal is to maximize 2 over the length of w and turn it into a problem where we can solve this directly.

(STILL) SUPPORT VECTOR MACHINES

- $\max \frac{2}{\|w\|}$ while classifying everything correctly
 $y_i (w^T x_i + b) \geq 1 \quad \forall_i$
- $\min \frac{1}{2} \|w\|^2$ quadratic programming
- $\max W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$
 $\text{s.t. } \alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

左下角的 \max 是 maximize 的意思

最後一行的 s.t. 應該是 subject to (the constraint) 的意思, 文中用過

6. >> Okay. So, we're still talking about support vector machines, although I haven't told you what support vector machines are yet, we're getting there, Michael. Bear with me. And what we got from our last discussion is that what we want to do somehow is maximize a particular equation, that is, 2 over the length of W . And as a reminder, W the parameters of our hyperplane. So somehow, we want to maximize that equation, subject to the constraints that we still classify everything correctly. Okay, so we want to maximize 2 over the length of W while classifying everything correctly. But, while classifying everything correctly is not a very mathematically satisfying expression, but it turns out we can turn that into a mathematically satisfying expression (即它下面那個式子 $y_i (w^T x_i + b) \geq 1$, 該式中, y_i 應該是指 x_i 這個點實際上應該為+還是-, 而不是前面那樣的 $y_i = w^T x_i + b$). And let me show you how to do that. So here's a simple equation. While classifying everything correctly turns out to be the same as, and I'm just going to write, I'm going to write it out for you, Michael, and see if you can, you can guess why this works. So, what I've written here is YI times W transpose XI plus B greater than or equal to 1 for all I . That is, for all of our training data examples. So why does this work? >> Well, what we really want is that the, that linear classifier, $WTXI$ plus B , is greater than or equal to 1 for the positive examples, and less than or equal to negative one for the negative examples. But you cleverly multiply it by the label on the lefthand side, which does exactly that. If YI is 1, it leaves it untouched. And if YI is negative 1, it flips everything around. So that we're really talking about less than or equal to minus 1. That's, that's very clever. >> It is very clever, and I'm going to pretend that I came up with that idea myself. So, it turns out that trying to solve this particular problem, maximizing 2 over W , while satisfying that constraint, is a little painful to do. But that we can solve an equivalent problem, which turns out to be much easier to do, and that is this problem. That is, rather than trying to maximize 2 over W , we can instead try to minimize 1/2 times W squared (並滿足那個 $y_i (w^T x_i + b) \geq 1$ 的約束. 注意此時 W 是 α 的函數, 我們要求的是 α , 它們使得 W 最大). Now can you see that those will always have the same answer? >> Yes, so, well, not the same answer, but it will be minimum, the point that maximizes one will minimize the other. Yeah, because the, we took the reciprocal. As long as we're talking about positive things. And since these are lengths, they'll be positive. Taking the

reciprocal exactly, you know, changes the direction, of what the answer is. And the squaring is, is, makes it monotone. It doesn't, it doesn't, it magnifies it but it doesn't change the ordering of things. So yeah. That, that, that seems fine. I don't why that's any easier, but it seems the same. >> Well, do you want to know why it's easier? Cause I'll tell you. >> Please. >> This is easier because when you have an optimization problem of this form, something like minimizing a W squared, subject to a bunch of constraints, that is called a quadratic programming problem. And people know how to solve quadratic programming problems in relatively straightforward ways. >> Awesome. >> Now, what else is nice about that is a couple of things. One is, it turns out that these always have a solution, and in fact have a unique solution. Now I am not going to tell you how to solve quadratic programming problems because I don't know how to do it other than to call it up in the MATLAB. But there's a whole set of classes out there, where they teach you how to do quadratic programming. We could take an aside, I could learn all about quadratic programming, and then we could talk about it for two hours. But it's really beside the point. The important thing is, that we have defined a specific optimization problem, and that there are known techniques that come from linear algebra that tell us how to solve them. And we can just plug and play and go. Okay? >> Okay, fair enough. >> Okay, fair enough. So, in particular, it turns out that we can transform, again, this particular quadratic programming problem into a different quadratic programming problem. Or actually, truthfully, into the normal form for a quadratic programming problem, that has the following form. So here's what this equation tells you, Michael. We have basically started out by trying to maximize the margin. And that's the same thing as trying to maximize 2 over the length of W , I think I convinced you of, subject to a particular set of constraints, which are how we codify that we want to classify every data point correctly in the training set. We've argued that that's equivalent to minimizing $1/2$ times the length of W squared, subject to the same constraints. And then notice, because we happen to know this, that you can convert that into a quadratic programming problem, which we know how to solve. And it turns out that quadratic programming problem has a very particular form. Rather than try to minimize $1/2$ of W squared, we can try to maximize another function that has a different set of parameters, which I'll call α . And that equation has the following form. It's the sum over all of the data points I , indexed by I , of this new set of parameters α , minus $1/2$ times, for every pair of examples, the product of their α s, their labels, and their values, subject to a different set of constraints. Namely that all of the α s are non-negative, and that the sum of the product of the α s, and the labels that go along with them, are equal to zero. >> Holy cow. >> Now, it's so obvious how you get from one step to the other I'm not going to bother to explain it to you. But instead tell you to go read a quadratic programming book. What I really need you to believe, though, mainly because I'm asserting it, is that these are equivalent. So if you buy up to the point that we are trying to maximize the margin, and that is the same thing as maximizing 2 over the length of W , and you buy that it's the same as trying to minimize $1/2$ times W squared, then you just have to take a leap of faith here that, if we instead maximize this other equation, it turns out that we are solving the same problem. And that we know how to do it using quadratic programming. Or other people know how to do it and they've written code for us. Okay? >> All right. >> All right, so trust me on this. This is what it is that we want to solve. Now, it turns out that we can run little programs to solve this, and you end up with answers. But what's really interesting is what this equation actually tells us about what we're trying to do. So let me just show you. This'll be just, talk a little bit about the properties of this equation, and the property of the solutions to this equation for a second. Okay? >> hm. >> Okay. So let me move a few things around so that we can look at it

(STILL) SUPPORT VECTOR MACHINES

$$\max W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

st. $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

- $w = \sum_i \alpha_i y_i x_i, b$
- α_i mostly 0 \Rightarrow few x_i matter

7. >> Okay. So we've done a little bit of moving, moving stuff around, and kept the same equation of before. Remember, our goal is to use quadratic programming to maximize this equation. So let me talk a little bit about the properties of the solution for this equation. So here's the first one. It turns out that once you find the alphas that maximize this equation, you can actually recover the w (即 $w^T x + b = 0$ 中的 w), which was the whole point of this exercise in the first place >> that's the little w , not the big W (要 maximize 的那個 $W(\alpha)$) >> That's the little W , not the big W . That's right, okay? >> Neat. >> Yeah, that is kind of neat. So it's really easy to do. And of course once you know w it's easy to recover b (由下一句知, 從 $w^T x + b = +1$ 中算得 b). You just find the value of X , you stick it into W , you, that you know is equal to plus 1, and then poof, you, you can find out B . So you can recover W directly from this and you can recover B from it in sort of the obvious way. But here are some other properties that are a little bit, little bit more interesting for you. So I want you to pay attention to two things. One I am just going to have to tell you, and the other I want you to think about. So here's the one that I'm going to tell you. It turns out, okay, that alpha, each of those α are mostly zero, usually. So if I told you that in the solution to this, most of the alphas that you come back are going to be zero, what does that tell you about W ? >> So W is the sum of the data points times their labels times α . And if the α is zero, that the corresponding data point isn't really going to come into play in the definition of W at all. So a bunch of the data just don't really factor into W . >> That is exactly right. So basically, some of the vectors matter for finding the solution to this, and some do not. So it turns out, each of those points are vectors. But you can find all of the support that you need for finding the optimal W in just using a few of those vectors. >> The non-zero alphas. >> Yeah, well the ones with non-zero alphas. So you basically built a machine that only needs a few support vectors. >> Oh. So the, the data points for which the corresponding alpha is non-zero, those are the support vectors? >> Yes, those are the ones that provide all the support for W . So knowing that W is the sum over a lot of these different data points, and their labels, and the corresponding alphas, and that most of those are zeroes, that implies, that only a few of the X 's matter. Now Michael, let me let me do a quick quiz. >> Yea, yea!

(STILL) SUPPORT VECTOR MACHINES

$$\max W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

st. $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

- $w = \sum \alpha_i y_i x_i, b$
- α_i mostly 0 \Rightarrow few x_i matter



8. >> Okay Michael, I've drawn a little teensy tiny graph on the screen. Can you see it? >> Yes. >> Okay, and some points are positive, some points are negative. And let's just imagine, for the sake of argument, that the green line that I've drawn in between them is in fact the optimal separator. It probably isn't, but let's just pretend that I drew the right one. Now, I've just told you that a lot of the alphas are going to be zero, and some of them are not. So, thinking about everything that we've said, and thinking about what it means to build, this, what it means to build, an, an optimal decision boundary, maximizing a margin. I want you to point to one of the positive examples that almost certainly is not going to have a non-zero alpha, and one of the minus examples that almost certainly is not going to have a non-zero alpha. That is, are not a part of the support vectors. >> You want something that does not have a zero? >> Don't confuse me. >> Do you want [LAUGH] something that's, that, that has a zero alpha or a non-zero alpha? >> I want something that has a zero alpha. >> Got it. >> That is, in some sense, doesn't matter. Okay, go.

9. Alright Michael. >> So. >> You think you got an answer? >> Yeah. It's interesting. So I guess it really does make some intuitive sense that the line is really, really nailed down by the points close to it. And the points that are far away from it, really don't have any influence. So I would put non zero, sorry, I would put zero alphas on the lower left hand minus and the, one of the upper pluses. >> Yes, and, you know, I haven't actually worked out the answer here. But, both of these pluses probably don't matter. Certainly, this one doesn't. certainly, this minus doesn't matter. Maybe this one doesn't matter, either. But the point that she raises, is exactly right. The, the points that are far away from the decision boundary, and can't be used to define the contours of that decision boundary. Don't matter, whether they're plus or minus. Does that make sense? >> Yeah, cool. >> Does this remind you of anything? >> Nearest neighbors? >> That's almost always the answer. Why does it remind you of nearest neighbors? >> because only the local points matter? >> Oh, that's a good answer. I was going to have a different answer. Know what my answer was? >> What? >> It's like kNN except that you already done the work of figuring out which points actually matter. So you don't have to keep all of them. You can throw away some of them. >> Oh, I see. So it doesn't just take the nearest ones, it actually does this complicated quadratic program to figure out which ones are actually going to contribute. >> Right, so it's just another way of thinking about instance-based learning, except that rather than being completely lazy, you put a lot, some energy into figuring out which points you could actually stand to throw away. >> Interesting. >> Okay. Yeah, I think that's kind of interesting. I think it's kind of cool. So good. So you got that.

(STILL) SUPPORT VECTOR MACHINES

$$\max W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{x_i^T x_j}_{\text{similarity}}$$

st. $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

- $W = \sum_i \alpha_i y_i x_i, b$
- α_i mostly 0 \Rightarrow few x_i matter

Well, let me show you one more thing, Michael. Alright, so you got this notion of there being very few of the, the support vectors that you need, but I want to point out something very important about some of the parameters in this equation. So we just got through talking about the alphas, right? Basically the alphas say, pay attention to this data point or not. But if you look carefully at this equation, the only place where the x s come into play with one another is here. So Michael, generally speaking, given a couple of vectors, what does $x_i^T x_j$ actually mean? >> It's the dot product. >> Right, and what is the dot product? >> It's like the projection of one of those onto the other right? >> Yeah, and that ends up giving you what? >> A number. >> Yes. Does that number kind of represent anything? And if you say the dot product, I will climb through the screen and kill you. >> What about the length of the projection. >> Right, and what does that kind of represent for you? >> Well, I guess in particular if the x 's are, well if there are five going to each other than it's going to be zero. But if they kind of point in the same direction, they're going to be, it's going to be a large value, and if they put in opposite directions it's going to be a negative value. So it's sort of kind of indicating how much they're pointing in the same direction. So, I guess it could be a measure of their similarity. >> Right, I think that is, that is exactly right. This is the kind of a notion of similarity. So if you look at this equation, what it basically says. Find all pairs of points. Figure out which ones matter for, for defining your decision boundary. And then think about how they relate to one another in terms of their output labels. With respect to how similar they are to one another.

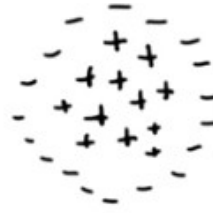
SVMs : LINEARLY MARRIED



10. >> Okay, Michael. So, I've drawn a little thing on the screen for you. because want to illustrate a problem. So, you see this little graph that I have? >> It looks just like all the other graphs you've drawn so far. >> More or less. I think there are fewer points, but I think you're right. It more or less looks like the same one. And, we found the line that is linearly separating the two clouds of points. >> [LAUGH] >> And you agree, that's a technical term. And you agree that it linearly separates them. And you're even willing to accept, for the purposes of this discussion, that is in fact the line that maximizes the margin. >> Sure. >> Even if it's not. Okay, cool. So, this is easy right? So, what are you going to do, Michael, if I take this now and I add one more point (即那個紅色的-)? And here is the point that I'm going to add. >> Hmm. >> It's another minus. >> I can think of two things to do. One is I can put a vertical line through it, and that makes things nearly separable again. >> That's usually not allowed. >> And the other one is, since it's a different color I feel like I could just erase it. >> No. >> No, all right. >> That's not acceptable. I control the pen. . >> I mean in some sense the margin is now negative, right. because this, this point, there's going to be no way of slicing this up so that all the negatives are on one side and all the positives on the other. >> That's true. >> It's not linearly separable. >> Okay. Can you think of some clever way to fix it? >> Well, again, I mean, I feel like, I mean I was making a joke before. But maybe a reasonable thing to do is to have some kind of, you know, I'm allowed to delete some number of points thing. Or, or, find the line that linearly separates the positives and the negatives, while at the same time, minimizing the number of things that are on the wrong side. >> Right. So you wouldn't be linearly summarizing, separating them. But you'd be finding a line that make sort of the minimal set of errors while also maximizing the margin, if you kind of were allowed to flip a few points from positive to negative or negative to positive. I think that's what you said. >> Yeah. And then you need to have kind of new knob now to trade off those two things. >> Right, and it turns out you can do that. We're not going to talk about it. Instead we're going to make a homework assignment about it then, and let the students think about it a little bit.

SVMs : LINEARLY MARRIED

$$W(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \textcircled{x_i^T x_j}$$



$$\Phi(q) = \langle q_1^2, q_2^2, \sqrt{2} q_1 q_2 \rangle$$

But I think even that little clever thing that you came up with, even though it is used, won't work in another case that I'm thinking of. So let me draw that case for you. Okay, Michael, how does that look? >> Well, the good news is it's not exactly like all the other graphs you've drawn but it is, it is very similar to it. There's going to be a linear separator that falls between that bottom plus and the line of minuses. And there's a maximum margin one. So its, this seems all within the bounds of what we've been talking about. << That's true. You're very smart. What if I add just a few more points? >> Okay, that doesn't seem like just a few. [LAUGH] Wait, so I guess [this is different from the other example because where before, as before, it looked like there was maybe like an outlier or an intruder. Now, it's like there's a ring around the whole thing. Oh, is that why they're linearly married?](#) >> Yes. >> Ha! All right. So now, you know, you can draw lines all day long, and it's just not going to slice things up. >> That's right. [So now we have to come up with some clever way of managing to make this work,](#) or we're going to have to throw away support vector machines altogether. And I like support vector machines, so I want to avoid doing that. So here's the little trick we're going to do. I am going to change the data points without changing the data points. >> Ouu. >> That's going to be a neat trick. >> That seems possible and impossible. >> Here's what I'm going to do, Michael. I am going to define a function, okay. Here's the function. [I'm going to create a little function here. And this function's going to take a data point,](#) okay. And because we've been using too many X's and I's and Y's and J's, [I'm simply going to call it q \(見下圖, 就知道 q 是甚麼意思了\),](#) okay. Now, [q is one of the points that are in,](#) it's in the same dimension as these other points. So in this case, it's two dimensions in a plane, okay. [And I'm going to transform that particular point q into another kind of point.](#) But I'm going to do it in another way that doesn't require cheating, okay. You ready? >> Sure, I'm perplexed but okay >> Okay. So what is Q? Q is in the Y, is in the plane. [So that means it has two different components, q1 and q2. And I am going to produce from those two components, q1 and q2, a triple.](#) So I am going to put that point into three dimensions now. And the dimensions are going to look like this. How does that look, Michael? >> Strange. So you took, so q/Q is a two dimensional point. So it's got, Q1 and Q2 are its two components. >> Yup. >> And you're saying, you're going to take the first component, make a new vector where the first component of that is squared, take the second component, make a new vector where the sec, that value is the second component squared. And now, just because apparently it wasn't weird enough, you're going to throw in a root 2 times the product of those two as the third dimension. Okay. >> That's right. >> You're, you know you have an interesting sense of style. >> I do. I kind of like it. Now let me point out something for you. One is [I haven't actually added any new information, in the sense that I'm still only using q1 and q2.](#) Yeah, I threw a constant in there, and I'm multiplying by one another, but at the end of the day, I haven't really done much. It's not like I've thrown in some boolean variable that gives you some extra information. All I've done is taken Q1 and Q2 and

multiplied them together, or against one another just because, why not? >> Okay. >> Okay? All right. Now, **Why did I do this? I did this because it's going to turn out to provide a cute little trick. And in order to see that cute little trick we need to return to our quadratic programming problem.** So let me remind you what that equation was. All right, Michael. So I've written up the equation for you, as I promised I would remind you, of the quadratic program that we're trying to solve. I didn't write down all of the constraints and everything. I'm hoping that you remember them. And I wrote it up there for a reason. And the reason I wrote it up, is because you'll recall, just not too long ago, I asked you to talk about XI and XJ, and what it looks like in this equation. And what I think we agreed to, or at least I know I agreed to, is that we can think about $x_i^T x_j$ as capturing some notion of similarity. So, it turns out then, if we sort of buy that idea of similarity, that really what matters most in solving this quadratic problem, and ultimately solving our optimization problem, is being able to constantly do these transpose operations. So, let's ask the question that, if I have this new function, fee or fi or whatever the right pronunciation is, what would happen if I took two data points, and I did the transpose or the dot product between them. What would I get? So, let me just write that out. Or, I don't know, you can try telling me if you want to. So let's make that, let's make that let's make that simple, Michael. And in fact, let's make it so simple we can make it into a quiz.

SVMs : LINEARLY MARRIED

$W(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j x_i^T x_j$

Quiz!

$x^T y \leadsto$

$\Phi(x)^T \Phi(y) =$

$\langle x_1^2, x_2^2, \sqrt{2} x_1 x_2 \rangle^T \langle y_1^2, y_2^2, \sqrt{2} y_1 y_2 \rangle = (x_1 y_1 + x_2 y_2)^2$

$\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2} x_1 x_2 \rangle$

$(x^T y)^2$

$x_1^2 y_1^2 + 2 x_1 x_2 x_1 y_2 + x_2^2 y_2^2$

Similarity

kernel trick

11. Okay, Michael. Here's a problem I want you to solve. **Let's imagine we have two points. I'm going to call them x and y, just so I can confuse you with notation. And they are both two dimensional points. So they're in a plane. And they have components x1 and x2 and components y1 and y2.** Okay? >> Sure. >> And **rather than computing the $x^T y$, their dot product, I want to compute the dot product of those two points, but passed through this function phi, or phi.** You got it? >> Yeah. Okay, so. X transpose Y, gets turned into phi X transpose phi Y. What's the answer? What's the output? In terms of X1, X2, Y1, Y2. Go.

12. >> All right Michael, you got the answer? >> I'm still carrying some squareds. >> You want to talk it through? >> Okay. Sure. So, x is really x1x2 and y is really y1y2 and phi x is now this crazy triple x so I, so I wrote x1 squared x2 squared, square root of 2x 1x 2. >> Yeah. >> That's the vector

that we get for $\phi(x)$. Then for $\phi(y)$, I get y_1 , it seems like it would be helpful to see this. >> You want me to write it down? >> Sure. >> Okay. So that turns out to be the same as what did you say? x_1^2 , x_2^2 . >> $\sqrt{2}$, x_1 , x_2 . >> $\sqrt{2}$, x_1 , x_2 . Okay. >> And then the y vector gets transformed to the same thing, except for with y s, y_1^2 , y_2^2 , $\sqrt{2}$, y_1 , y_2 . >> Okay. >> So, then, the, the dot product is just, the, the products of the corresponding components summed up. So x_1^2 , y_1^2 plus, >> Okay >> x_2^2 , y_2^2 , >> Mm-hm x_1 , x_2 , y_1 , y_2 >> That's right. So here's a question for you, Michael. Does that look familiar? Based on your years of thinking about algebra. >> Oh, thanks for writing it that way! I see. We can, is this right? So it's, it's like, we can factor this. >> Mm-hm. >> It's like x_1 plus x_2 times y_1 plus y_2 . >> Yeah, that's right. Wait, is that right? >> No it's not right. $x_1 y_1$ >> Mm-hm. >> Plus $x_2 y_2$. Whole thing's squared. >> Right. So, let's write that down. So if you factor it out, you're right, this is exactly equal to, $x_1 y_1$ plus $x_2 y_2$. Squared. And what's an even simpler way of writing that? >> I see. $x^T y$ looks like a dot product itself. It looks like the dot product of x and y . Oh, it's $x^T y$ on the inside, and then we square it on the outside. >> That's exactly right. So now do you see why there was method to my madness when I created the ϕ function? >> Not yet. So you're saying, if we're just dealing with dot products, now I'm still a little confused. So, so it is the case, that you define these in an interesting way, so that the dot product became the square of the old dot product. >> Right, so now let me make two observations. Okay. Here's observation one. What's $x^T y$? I mean geometrically, what does that represent? >> The length of the projection of y onto x . >> No I mean geometrically, like go all the way back to geometry, third grade. >> We didn't do transposes in third grade. [LAUGH] >> I know we didn't do transposes, but you did equations like this, or at least later you learned they were equations like this. Here, pretend you're in third grade, and I said talk to me about geometry. What kind of words would you use? >> Oh, what's in geometry? Triangles. Circles. >> Yeah. Circles! Did you say circles? >> Sure, but only because I thought you might have wanted me to. >> Did you say circles? >> That's a really ugly looking circle, sure. >> Sure. This is basically a particular form of the equation for a circle, which means that we've gone from thinking about the linear relationship between x_i and x_j or your data points and we've now turned it into something that looks a lot more like a circle. >> Interesting. >> So if and this is my second point, Michael. If you believe me in the beginning where we notice that $x_i^T x_j$ is really about similarity. It's really about why it is you would say two points are close to one another or far apart from one another. By coming into this transformation over here, where we basically represented the equation for a circle, we have now replaced our notion of similarity from being this very simple projection to being the notion of similarity is whether you fall in or out of a circle. >> So more sort of about the distance as opposed to what direction you're pointing. >> Right, and both of them are fine because both of them represent some notion of similarity, some notion of distance in some space. In the original case, we're talking about two points that are lined up together. And over here together with this particular equation we represented whether they're inside the radius of a circle or outside the radius of a circle. Now this particular example assumes that the circle is centered at the origin and so on and so forth, but the idea I want you to see is that we could transform all of our data so that we separated points from within one circle to points that are outside of the circle. And then, if we do that, projecting from two dimensions here into three dimensions, we've basically taken all of the pluses and moved them up (out of the page) and all of the minuses and moved them back (into the page), and now we can separate them with the hyperplane. >> Without knowing which ones are the pluses and which ones are the minuses, of course. >> Of course. >> Because, I see, because they are the ones that were closer to the origin. >> Right. >> So they get raised up less. Wow. Okay. So using that third dimension. >> Right. Now, this is a cute trick, right? I can basically take my data, and I transform it into a higher dimensional space, where suddenly I'm now able to separate it linearly. That's very cute, but I chose this particular form for a reason. Can you guess why? >> Because it fits the circle pattern that you wanted. >> But there are lots of different ways we could have fit the circle pattern. I chose this particular form because not only does it fit the circle pattern, but it

doesn't require that I do this particular transformation. Rather than taking all of my data points and projecting them up into three dimensions directly, I can instead still simply compute the dot product and now I take that answer and I square it (因為 $\Phi(x)^T \Phi(y)$ 得到的也是這個結果). >> So you're saying in this formulation of the quadratic program that you have there in terms of capital W if you write code to do that, each time in the code you want to compute that $x_i^T x_j$, if you just squared it right before you actually used it, it would be as if you projected it into this third dimension and found a plane? >> Yes, that's exactly right. >> That's crazy. >> It's so crazy it has a name. And that is the kernel trick. So, again if we really push on this notion of similarity. What we're really saying is we care about maximizing some function that depends highly upon how different data points are alike, or how they are different. And simply by writing it this particular way, all we're saying is, you know what, we think the inner product is how we should define similarity. But instead, we could use a different function altogether, Φ or more nicely represented as $(x^T y)^2$, and say, that's our notion of similarity. And we can substitute it accordingly. >> So we never really used phi. >> We never used Φ . We're able to avoid all of that by coming up with a clever representation of similarity. That just so happened to represent something, or could represent something in a higher dimensional space. >> So is it important that such a phi exists? Or is it just the case that we can, you know, we could, we could throw in a cubed, we could throw in a fourth, we could do a square root and a log, like can we do anything we want there in that $x_i^T x_j$? Or are we constrained to only use things that somewhere out there, there is a way of representing it as a regular dot product? >> Well, that is an interesting question. The answer is (for the kernel) you can't just use anything, but in practice it turns out you can use almost anything. And the other answer to your question is, it turns out for any function that you use, there is some transformation into some dimensional space, higher dimensional space, that is equivalent. >> Whoa. >> Now, it may turn out that you need an infinite number of dimensions to represent it. But there is some way of transform, transforming your points into higher dimensional space that happens to represent this kernel, or whatever kernel you choose to use. >> So, which part is the kernel? >> So, the kernel is the function itself. So, in fact, let me, let me, let me clean up this screen a little bit. And, and see if we can make this a little bit more precise and easier to understand.

$$W(\omega) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{(x_i^T x_j)}_{\text{kernel}}$$

SVMs : LINEARLY MARRIED

$$W(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underline{K(x_i, x_j)} \quad \begin{array}{l} \rightarrow \text{similarity} \\ \rightarrow \text{domain} \\ \quad \underline{\text{knowledge}} \end{array}$$

$$K = (x^T y)^2$$

$$K = x^T y$$

$$K = (x^T y + c)^p$$

$$K = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

$$K = \tanh(\alpha x^T y + \theta)$$

Mercer Condition

13. >> Okay, so I've cleaned up the screen a little bit, Michael to, to make this a little bit clearer. Now, let's look at this $x_i^T x_j$. And I'm now going to replace it with something. So, I've just replaced it with a function, which I'm going to call a kernel. Which takes x_i and x_j as parameters, and will return some number. And again, as we talked about before, we think of the $x_i^T x_j$, as some notion of similarity, and so this kernel function is our representation, still, of similarity. Another way of thinking about that, by the way, is that this is the mechanism by which we inject domain knowledge into the support vector machine learning algorithm. >> Just like we were injecting domain knowledge when we were thinking about k-nearest neighbors. >> Yes, everything has domain knowledge and everything ultimately comes back to k-nearest neighbors. I don't know why and I don't know how, but it always seems to. >> So the k in k-nearest neighbors, and the k in kernel really stand for knowledge. >> Oh, wow, that's pretty good. We should write a paper with that title. >> [LAUGH] So the room neatness here, the neatness here two fold. One is, you can create these kernels and these kernels have arbitrary relationships to one another. so, what you're really doing is, projecting into some higher dimensional space, where, in that higher dimensional space, your points are in fact, linearly separable. But, the second bit is, because you're using this kernel function to represent your domain knowledge, you don't actually have to do the computation of transforming the points into this higher dimensional space. I mean, in fact if you think about it, with the last kernel that we used, computationally, there was actually no more work to be done. Before we were doing x transpose y , and now we're still doing x transpose y , except we're then squaring it. So that's just a constant bit more work. Right? >> So is, and that is a kernel and another kernel is X transpose Y ? >> Yes, that's something I would call a kernel. >> And the other kernel we talked about was just X transpose Y by itself. That's, that's a kernel too, isn't it? >> Oh no, no. That's right. That's right. That's absolutely right. So, that's a different kernel, you're absolutely right. Just X transpose Y . Is another kernel. Actually we can write a general form of both of these. And as a very typical kernel, it's the polynomial kernel where you have x transpose y plus some constant, let's call it c , raised to some power p . And as you can see, both of those earlier kernels are, in fact, just a special case of this. >> Hm, and I would, yeah, okay, good. >> And that should look familiar. >> It reminds me of the regression lecture. >> Exactly, where we were doing polynomial

regression. So now, rather than doing polynomial regression the way we were thinking about it before, we use a polynomial kernel and that will allow us to represent polynomial functions. And there're lots of other kernels you can come up with, Michael. So. Here's just a couple. I will just sort of leave em. Leave em up to you, to think about. And there's, there's tons of them. So, here's one that I, I happen to like. So, that's a sort of, radial basis kernel. Does that look familiar to you? >> Well, to me, make sure I understand that it's doing the right thing. So if x and y are really close to each other, then it's like, e to the minus zero over something which is like e to the zero, which is like one. So there's similarities like one if they're on top of each other. If they're very far apart, then it's like their distance is something very big divided by something e to the minus something very big is very close to zero. So it does have that kind of property kind of like the sigmoid where it, it transitions between zero and one but it's not exactly the same shape as that. >> Right in fact it's symmetric, that the square of the of the distance between you in making an actual distance, makes it always a positive value there. Or at least a non-negative value there. And so it becomes symmetric, so it looks a lot more like a, like a Gaussian with some kind of width which is represented by sigma. And there are tons and tons of these. Actually if you wanted to get something that looked like a sigmoid, here's one. Where alpha's different from the other alphas, but I couldn't think of a different Greek letter. And this function gives you something that looks a lot more like a sigmoid. And there're tons and tons of these you can come up with. And there's lots of, been a lot of research over the years on what makes a good kernel function. The most important thing here, I think, is that it really captures your, your domain knowledge. It really captures your notion of similarity. You might notice, Michael, that since it's just an arbitrary function that returns a number, it means that X and Y or the, the different data points you have, don't have to actually be points in a numerical space. They could be discrete variables. They could describe whether you're male or female. As long as you have some notion of similarity to play around with, that you can define, that returns a number, then it doesn't matter. It will always work. >> So can you do things like, I don't know, strings or graphs or images? >> Absolutely. You could think about two strings. How are two strings similar? Maybe they're, they're similar if their edit distance is small. The number of transformations that you have to give in order to transform one string to another. If there are few of those, then they're very similar. If there are a lot of those then they're very dissimilar. You could talk about words like cat and lion, and decide those are more similar than cat and mosquito, for example. >> because they, because of the ears? >> Yeah. Mainly because of the ears. >> All right. But then I, then I think I understand. >> Okay. Good. So you might be curious, Michael, whether there are any bad kernel functions. There is actually an answer to that. While it's not clear whether there are any bad kernel functions, it is the case that in order for all the math to go through, there is a specific technical requirement of a kernel function. It has a name. And it's the Mercer Condition. Have you ever heard of the Mercer Condition? >> I've heard the word. I actually used to live near Mercer County in New Jersey. >> You did? >> Yeah. >> Oh. So then I guess it's the condition of living near where Michael used to live. Now, so the Mercer condition is a very technical thing we'll talk about this again, a little bit in the homework assignment. But for your intuition in the meantime, it basically means it acts like a distance, or it acts like a similarity. It's not an arbitrary thing that doesn't relate the various points together. Being positive is something definite in in this context means it's a well behaved distance function. >> Gotcha.

SVMs

- margins \sim generalization \neq overfitting
- big is better
- optimization problem for finding max margins: QPs
- support vectors
- Kernel trick $x^T y$ $K(x, y)$ \leftarrow domain knowledge

14. Okay, Michael, [so that gets us to the end of Support Vector Machines. So, let's recap. What have we learned?](#) >> Well, we learned that support vector machine is not a command or a political statement. We talked about how a margin is a, is a useful concept in trying to understand how well a linear classifier might generalize. >> Okay, I like that. Lemme write that down. Margins are important. So we learned about margins In particular their relation to generalization and overfitting. >> Okay. >> In particular, we, we would like, given the choice, to find a linear separator that has the largest margin. >> Right, right. >> So maximizing margins. >> Right. At least when it comes to margins, bigger is better. >> Then we talked about how you could actually find. A linear separator that has maximum margin. I think we turned it into a quadratic program. >> Yes. We found an optimization problem for finding maximum margins and they turned out to be quadratic programming. >> And it was the dual of the quadratic program that showed us how, what the support vectors were. The support vectors were the points from the input data. That were necessary for defining that maximum margin separator. >> Right, right. So, we actually figured out what support vectors were. And then we tied all that in to instance based learning and other kind of ensemble methods. And so you could sort of think of support vector machines as being eager lazy learners. Or only as lazy as necessary to represent what you needed to represent. Because the, the classifier was based on just a subset of the data, or, or, or the raw data was being used for defining the classifier. >> Exactly right. Alright, so is there anything else? >> Oh. Oh. Right. Then there was the whole idea that, well, linear doesn't always seem like enough, but, we can project. Data into a higher dimension space and do the comparisons there. And that only made one little change to the algorithm. In particular, the dot product could be turned into some other similarity metric. And you called that the kernel trick. >> Right. >> love the kernel trick. And as you say, we took, basically, $X^T Y$. And generalized it to a generic similarity function k of x comma y . And that actually ended up representing all of our domain knowledge, which will come up again and again throughout this course. What are the levers we have For expressing domain knowledge. Okay so, anything else for writing down Michiel? >> I think you said that kernels had to satisfy the merciful condition. >> No, no, no, mercer condition. >> Oh, okay, well that makes more sense. >> [The mercer condition is interestingly quite merciful, because it's ok to use kernels that don't necessarily](#)

satisfy the mercer condition, often in practice it still works. Okay, well good, I think that is mostly it. So i guess we're done. >> Wait a second. Didn't you say that you were going to explain boosting to us? You kind of suggested maybe you'd do that during the boosting lecture. But then you said it would be during the SFM lecture. And now the SFM lecture's over and we still don't know why boosting doesn't over fit. >> Oh. That's a good point. That's a good point. I had forgotten about that. Thank you Micheal, this is why I keep you around, to remind me about stuff like this. Okay, let me take a moment then, just a moment. To see if we can tie together, over fitting and boosting about what we just learned about SVM's.

Back to Boosting



15. >> Alright, so back to boosting, Michael. So as you recall the little teaser I left you with last time, is that it appears that boosting does not always over-fit. And a little graph. That's true, but it doesn't seem to over-fit in the ways that we would normally expect it to over-fit. And in particular we'd see a, you know, an error line on training And what we expect to see is a testing line that would, you know, hue pretty closely and then start to get bad.

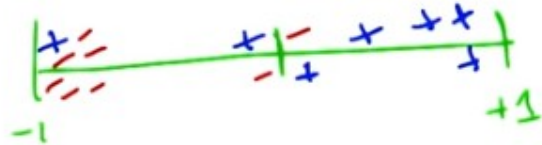
BACK TO BOOSTING



+ error

+ confidence

$$H(x) = \text{sgn} \left(\frac{\sum_b \alpha_b h_b(x)}{\sum_b \alpha_b} \right)$$



But what actually happens is that instead, [this little bit at the end where you get over fitting seems to instead. Just keep doing well. In fact, getting better and better and better. And I promised you an explanation for why that was.](#) So, given what we talked about with support vector machines, and what we spent most of our time thinking about, what do you think the answer is? Well I, I don't think I would have asked again if I, had a thought about it. But you mean You want me to connect it to support vector machines, somehow. Well the, the thing that was fighting over fitting in support vector machines, was trying to focus on maximum margin classifiers. >> Here, let me, let me try to explain to you why it is that you don't have this problem with With overfitting at least not in the, in the typical way as you keep applying it over and over again like you do with something like neural networks. And it really boils down to noticing that we've been ignoring some information. So, what we normally keep track of is error. So error on say a training set is just, you know, the The probability that you're going to come up with an incorrect answer or come up with an answer that disagrees with your training set. and that's a very natural thing to think about and it makes a lot of sense. But there's also something else that is actually captured inside of boosting and captured by a lot of learning algorithms we haven't been taking advantage of, and that's the notion of confidence. [So confidence is not just whether you got it right or wrong. It's how strongly you believe in a particular answer that you given.](#) Make sense? >> Yes, a lot of the algorithms we talked indirectly have something like that. So, [like in a nearest neighbor method, if you are doing five nearest neighbor and all five of the neighbor agree, that seems different than the case with vote one way and two vote the other.](#) >> Right. And in fact, that's a really good example. If you think of that in terms of regression Then you could say something like the variance, between them is sort of a stand in for confidence. Low variance means everyone agrees, high variance means, there's some major disagreement. Okay. So what does that mean in the boosting case? Well as you recall, [the final output of the boosted classifier](#) is given by a very simple formula. And here's the equation here that h of x is equal to the sign of the sum over all of the weak hypothesis that you've gotten of α times h . So the weighted average of all of the hypothesis, right? And you just simply, if it's positive you produce a plus one. And if it's it negative you produce a minus and if it's exactly zero

you don't know what to do so you just. Produces zero. Just throw up your hands. So I'm going to make a tiny change to this formula, Michael. Just, just for the purpose of sort of, explanation, that doesn't change the fundamental answer. And I'm just going to take exactly this equation as it is. And I'm going to divide it, by the weights that we use. Now what does that end up doing? >> Okay, so the weights. I'm getting a. There's Alphas in the SVM's too, so I'm getting a little confused. So that I'm. I think these Alphas all have to be non-negative. >> Right. >> But they kind of like this support vector values, in that there could be zero, if, if that hypothesis isn't come into play? >> Well, but they want in that case, the, the alpha is always set to be the natural log of something. >> Oh, oh, oh, and also these alphas are applied to hypothesis whereas the alphas in the, in the SVM settings were being applied to data points. >> That's right. So, unfortunately in machine learning, people in, invent things separately and re-use notation. Alpha's an easy Greek character to draw, so people use it all the time. But here, remember, alpha's the measure of how good a particular weak hypothesis was, and since it has to do better than chance, it works out that it will always be greater than zero. >> Gotcha, okay. So this, this normalization factor, this denominator doesn't, it's just a constant with respect to x, the input. So it won't actually change the answer. So it really is the same answer as we had before, just a different way of writing it. >> Right. And what it ends up doing like often is the case in these situations, is it normalizes the output. So it turns out that this value. Inside here (即 `sgn()` 的括號裡) is always going to be between -1 and +1. Okay? But otherwise it doesn't change anything about what we've been doing for boosting. So you might ask why did I go through the trouble of normalizing it between minus one and plus one? >> Why indeed? >> Well it's makes it easier for me to draw what I want to draw next. So, we know that the output of this little bit inside the sign function is always going to be between minus one and plus one. Let's imagine that I take some particular data point x and I pass it through this function, I'm going to get some value between minus one and plus one. And let's just say for the sake of the argument, it ends up here (+1 附近). Okay? >> Is that an x or a plus? >> That's a plus. >> Okay. So it's a positive example and it's near plus one (那些藍色的+和紅色的-表示實際(即正確)的結果; 而藍色的+或紅色的- 所在的數軸上的坐標值 表示我們算出的結果, 不一定正確). >> Right. >> So this would be something that the algorithm is getting correct. >> Yes, and it's not just getting it correct, but it is very confident in its correctness, because it gave it a very high value. By contrast there could have been another positive that ends up around here (0 附近). >> Hmm. >> So it gets it correct but it doesn't have a lot of confidence so to speak in its correct answer because it's very near to zero. So that's the difference between error and confidence. Because for example I could also have a plus value way over here (在-1 附近). So I am very, very confident in my very, very incorrect answer. >> Mm. >> So this is my daughter, for example. [LAUGH] She's very confident whether she's right or wrong. [LAUGH] Okay. And so now imagine there's lots of little points like this. And if you're doing well, you would expect that, you know, very, very often you're going to be correct. And so you end up shoving all the positives over here to the right, and all the negatives over here to the left. And it would be really nice if you were sort of confident in all of them. Okay, so does this make sense, Michael as a picture, >> Oh yeah. >> What, what might be going on? Absolutely.

BACK TO BOOSTING



+ error

+ confidence

$$H_{\text{final}}(x) = \text{sgn}\left(\frac{\sum \alpha_t h_t(x)}{\sum \alpha_t}\right)$$

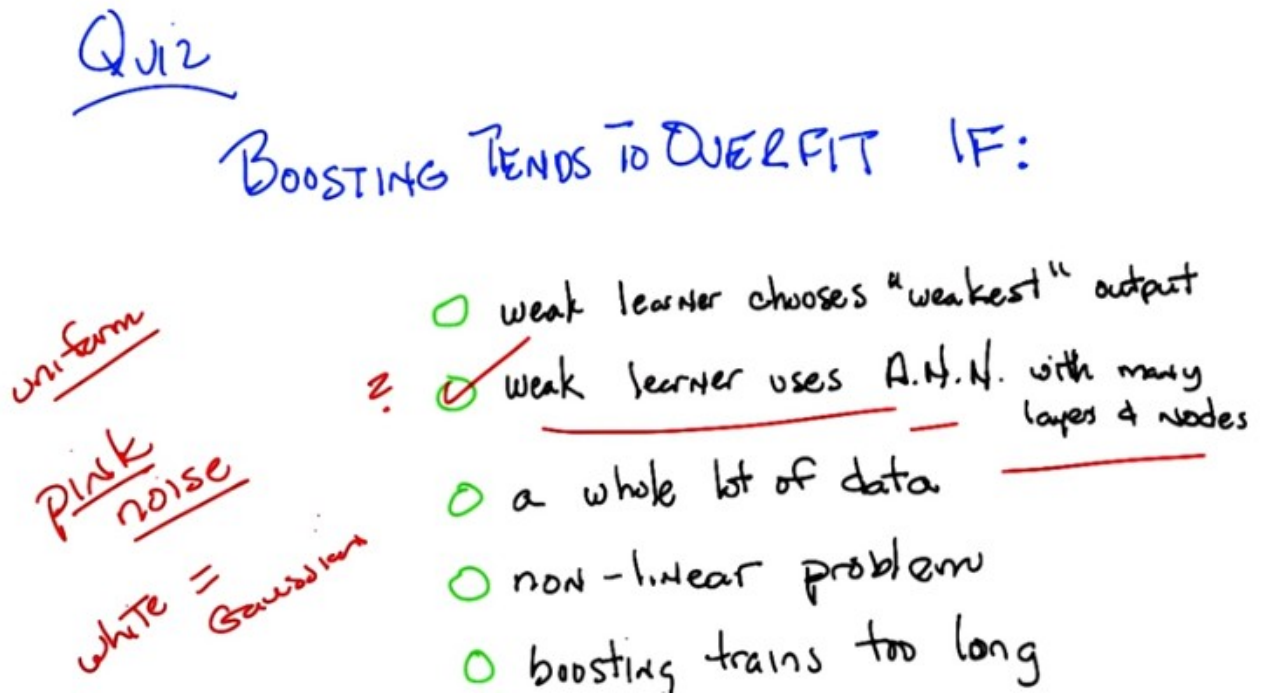
margin

-1

+1

>> Okay, good. So now I want you to imagine that we've been going through these, these training examples, and we've gotten very, very good training error. In fact, let's imagine that we have negative training error. I'm [LAUGH] >> Wow. >> In fact, let's imagine that we have no training error at all. So we, we label everything correctly. So then the picture would look just a little bit different. We're going to have all the pluses on one side, and all the minuses on the other. But we keep on training, we keep adding more and more weak learners into the mix. So here's what ends up happening in practice, right? What ends up happening in practice is, you have to do some kind of distribution on the hard examples. And the hard examples are going to be the one that are very near the boundary (0 附近). So as you add more and more of these weak learners what seems to happen in practice is that these pluses that are near the boundary and these minuses that are near the boundary just start moving farther and farther away from the boundary. So, this minus starts drifting and drifting and drifting until it's all the way over here, this minus starts drifting and drifting and drifting until it's all the way over here. And the same happens for the pluses. And as you keep going and you keep going, what ends up happening is that your error stays the same. It doesn't change at all, however your confidence keeps going up and up and up and up and up. Which has the effect, if you'll look at this little drawing over here of moving the pluses all around over here, so they're all in a bunch, and the minuses are on the other side. So what does that look like to you, Michael? >> This picture? >> Yeah. >> I mean that there's a, there's a big gap between the left most plus and the right most minus. Which, you know, in the context of this lecture reminds me of a margin. >> That's exactly right. Basically what ends up happening is that as you add more and more weak learners here the boosting algorithm ends up becoming more and more confident in its answers which it's getting correct. And therefore effectively ends up creating a bigger and bigger margin. And what do we know about large margins? >> Large margins tend to minimize over fitting (本文件最前面講過). >> That's exactly right. So it, counter intuitively, as we create more and more of these hypotheses, which you would think would make something more and more complicated, it turns out that you end up with something smoother, less likely to overfit and ultimately, less complicated. So the reason boosting tends to do well and tends to avoid over fitting even as you

add more and more learners is that you're increasing the margin. And there you go. And [if you look in the reading that we gave the students there's actually a detailed description about this in a proof.](#) >> Cool. >> Okay. So, there you go, Michael. Do you think, then, that boosting never overfits? >> [SOUND] Never seems like such a strong word. I mean, the story that you told says that it's going to try to separate those things out, but I guess I guess it doesn't have to be able to do that. I mean, it could be that for example all the weak learners are I dunno very unconfident very inconsistent. >> Hm. Okay, well you know, maybe, maybe it's worthwhile to take a little diversion here to take a five second quiz. >> I think it's worth the time. >> Done!



第二個選項中的 A.N.N.是 artificial neural network.

16. Okay Michael. So here's a quick quiz. [So we just tried to argue that boosting has this annoying habit of not always over fitting, but of course something can always over fit. Because otherwise we just do boosting and we're done, then neither of us would have jobs. And we don't want that to happen. So here's a little quiz to see if we can figure out the circumstances under which boosting might over fit, or tends to over fit. So here are five possibilities. Let me read them to you. Tell me if they actually make sense. So here's possibility number one, boosting will tend to over fit if The weak learner that it's boosting over, always chooses the weakest output that is, it, among all the hypothesis that it finds that do better than chance over the training with whatever given distribution. It always pick the one that is still nonetheless closest to chance, while still being better.](#) >> Well, why would it do that? >> Just to be difficult. >> Alright, and so you want to know, whether that makes it over, would [INAUDIBLE] make it over fit? [UNKNOWN] boosting overfit. >> Okay. Alright. >> The second one is weak learner actually ends up using...or the weak learner itself that boosting is using is in fact a neural network learner. And just for a little specificity, let's say this is a neural network that has many many layers and many many nodes. So, you know, it's a big powerful neural network, alright? the other option is... Boosting has a lot of data. So you're trying to learn, your training data is actually very, very, very large. You have lots and lots of examples. The fourth case, is that, the true underlying hypothesis, the true underlying concept, is in fact non linear. So you can't just draw a line. And then the

fifth case is that we let boosting train much too long. Whatever that means. Let's just say we let it train a lot. Not just a thousand iterations but a hundred billion iterations. >> Okay. >> Okay. Billions and billions of iterations. Okay. You got it? >> Yeah. >> Alright. Go.

17. >> All right, Michael. What's the answer? >> All right. Well, let me start off with what I think the answer isn't. So, the last one, boosting tends to overfit, if boosting trains too long. You just told me a story about that not being true. So I'm going to eliminate that one from consideration. Boosting training too long. >> Oh, nice to know you were listening. >> [LAUGH] Boosting training too long, seems like not a good reason for it to overfit. >> You're correct. >> All right. Boosting tends to overfit if it's a nonlinear problem. So, that doesn't seem right. I mean I guess, no, this one just doesn't seem right at all. Like I don't see why, why the problem being linear or nonlinear, has anything to do with overfitting., >> Okay. >> A whole lot of data is the opposite of what tends to cause overfitting. If there's lots of data then you'd think that it would actually do a pretty reasonable job of, you know, there's a lot to fit. There's a lot going on there. It's unlikely to overfit. >> Right, and in fact if a whole lot of data included all of the data, and you actually could get zero training error over it, then you know you have zero training zero generalization error. >> because it'll work on the testing data as well, because it's in there. >> Right. >> All right. Weak learner uses artificial neural network with many layers and nodes. So I'm guessing that you wanted me to think about that being something that, on its own, is prone to overfitting, because it's got a lot of parameters. >> Sure. >> So, if, and now we're doing boosting over that. So we fit a neural net, and then we fit another neural net, and we fit another neural net. And we're combining all the outputs together in the correct, weighted way. It's not obvious to me that that should be a good thing to do. I'm not sure it would overfit, but it seem like it sure could. >> OK, so you're, you're, so for now let's put a little question mark to it. You think that might be the right answer, but you want to think about it some more? >> Yeah let me, let me look at the first one. Weak learner chooses the weakest output. Well, I mean boosting is supposed to work as long as we have a weak learner. . And it doesn't matter if it chooses the weakest or the strongest. All that matters is it does significantly better than a half. So, like I feel like the only one, [the only one of these choices that is likely to be true is the second one.](#) >> And that is, in fact, correct. [So let me give you an example of when that would be correct. So let's imagine I have a big powerful new network that could represent any arbitrary functions. Okay, got lots of layers and lots of nodes. So, boosting calls it, and it perfectly fits the training data, but of course overfits. So then it returns, and it's got no error \(think of polynomial regression, 當 order 大時, 出現 overfit, 但 overfit 時 對於個 sample 數據點的 error 是很小的, 只是曲線形狀 變得不對\), which means all of the examples will have equal weight. And when you go through the loop again, you will just call the same learner, which will use the same neural network, and will return the same neural network. So every time you call the learner, you'll get zero training error, but you will just get the same neural network over and over and over again. And a weighted sum of the same function is just that function. So if it overfit, boosting will overfit.](#) >> Interesting. And not only will it overfit, but it'll just, it'll be stuck in a horrible loop of error. >> Right. [So that's why this is the sort of situation where you can imagine boosting a lower fit. If the underlying learners all overfit and you can never get them to stop overfitting, then there's really not much you can do.](#) >> Interesting. >> Now, I do want to have a little semantic argument with you for a moment, Michael. You used the word strongest at some point, when you were talking about using the weakest output. And I just want to point out that, that doesn't really mean anything. >> What do you mean, it doesn't mean anything? >> Well, so what's a strong, what would you call a strong learner? >> One that is far away from it. If a weak learner just has to do a little bit better than a half, it seems like a strong learner would be something that would be very close to being accurate. >> Right. Of course, on the other hand, if by that definition all strong learners are also weak learners. >> Sure. >> Because anything that does better than a half is still doing better than a half, which is all it requires to be a weak learner. >> Yeah, but that's kind of true of people too. Like a strong person is also a weak person. >> No. >> Well it depends how you define it. So, if you

say a weak person is someone who can at least lift their own arms, then strong people are also weak people in that they can lift their arms. >> Yes if you define it that way and if I define blue to be purple, then I can say blue is purple. But that's not how people define weak people. They define weak people, by saying they can't lift more than, not that they can lift at least as much. >> I see. So it's this piece of terminology that boosting uses that is in error, not me. >> That's one interpretation. It's not the one that I would use, but it's one interpretation. When you say something like a strong learner, I mean, it makes sense to use that kind of term, and sort of throw it around, and say, well, by a strong learner I mean someone who's, or a learner that's going to overfit, or is going to always do really well on the training data. But in kind of a technical definition it's very difficult to sort of pin down. So don't get too caught up what a strong learner means if you want to write a proof. Seems fair? >> Good point yeah, also, also that this whole notion that strong is sometimes defined as not weak. And it is not the case that if you have something that's not a weak learner that it's, then it's a strong learner. In fact, it's no learner at all. >> Exactly. So, a weak learner's just defined in a way that basically says, it gives me at least some information. Good. Let me just throw one more thing in here and then we can stop talking about this. [There's another, a couple of other cases where boosting tends to overfit. The one that matters the most, or comes up the most, is in the case of pink noise.](#) >> Did you say, peak noise? >> I said, pink noise. I even wrote it in red, which looks like pink. It's a strong pink as opposed to a weak pink. >> [LAUGH] >> I'm sorry. There's no way for that to be obvious from what we've talked about, but as a practical matter, pink noise tends to, cause boosting overfit. >> Okay, but this is not a term I'm familiar with unless you're critiquing the musical stylings of a particular performer. >> [LAUGH] No. Although I did recently see, see them in concert. But that's a whole other conversation. Okay, so pink noise just means uniform noise. >> I thought white noise was uniform noise. >> No, white noise is Gaussian noise. Okay, so [pink noise is uniform noise and white noise is Gaussian noise.](#) This is why, Michael, by the way, if you ever try to set up a studio or a cool stereo system in your house, you want a pink noise generator. So that it covers all the frequencies equally, not just the white noise. generated. >> Hm. >> [But boosting tends to overfit in those sorts of circumstances. And you can read more about it in the notes if you want to. But the one that I really want people to get is, that if you have an underlying weak learner that overfits, then it is difficult for boosting to overcome that. Because fundamentally you've already done all of your overfitting and it's, there's really not much for those things to do.](#) >> Okay. Got it? >> Got it. >> Excellent. It all ties back into margins, and it's all one big story, which I think is the lesson of all of machine learning.

SVMs

- margins ~ generalization & overfitting
- big is better
- optimization problem for finding max margins: QPs
- support vectors
- Kernel trick $x^T y$ $K(x, y)$ ← domain knowledge

18. >> All right Michael, so we're back. We have learned a bunch of stuff about SVMs and you forced me to live up to my promise to connect margins. Which is kind of the big cool thing about SVMs, the sort of fundamental theoretical construct that's underlying SVMs. You forced me to connect that back to boosting. So I think we've learned everything we need to learn about support vector machines for now. What do you think? >> If you say so, sounds good. >> All right. Well, I will talk to you next time, Michael. >> Thanks. Bye.