

Support Vector Machines

Lecture 10

Termeh Shafie

1

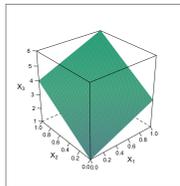
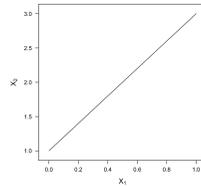
Support Vector Machine (SVM) is a supervised learning algorithm used to learn a **hyperplane** that can solve the **binary classification problem**

2

Hyperplanes

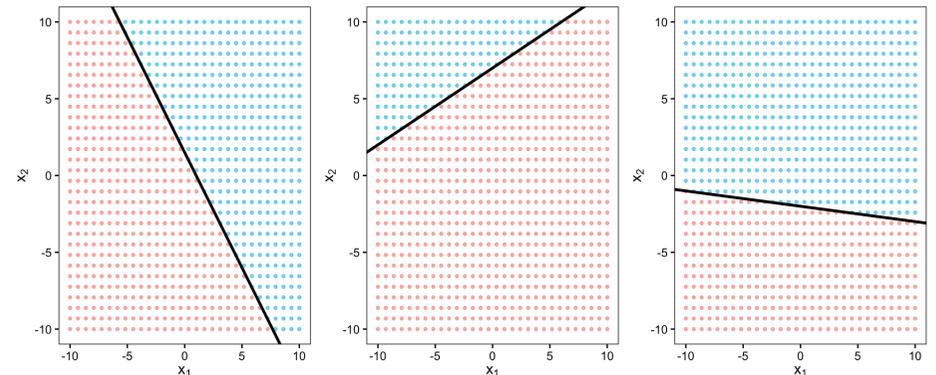
"a flat affine subspace" 😊

- **Flat:**
 - hyperplane is not curved, it increases/decreases constantly in each direction
- **Affine:**
 - the hyperplane doesn't need to pass through the origin
 - it can have an "offset" or be shifted (may have intercept)
- **Subspace:**
 - a subset of vectors in a larger vector space
 - in a d -dimensional space, a hyperplane has dimension $d - 1$
 - in 3D it is a **plane**
 - in 2D it is a **line**
 - in 1D it is a **point**



3

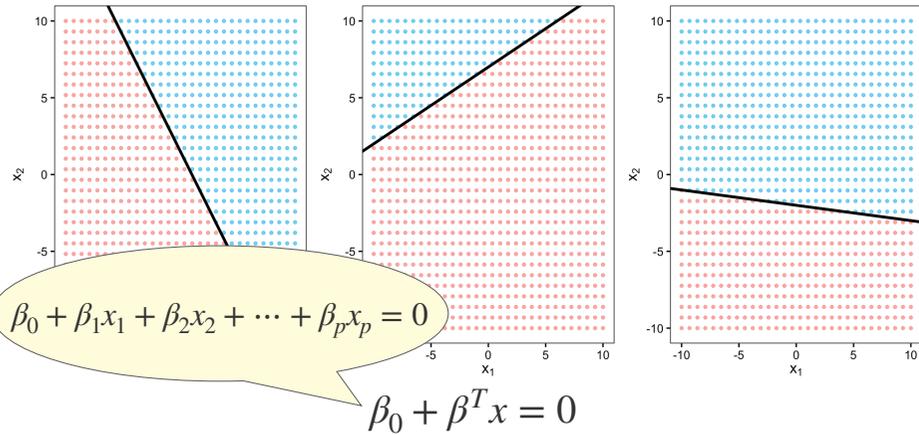
Hyperplanes Divide the Space in Half



$$\beta_0 + \beta^T x = 0$$

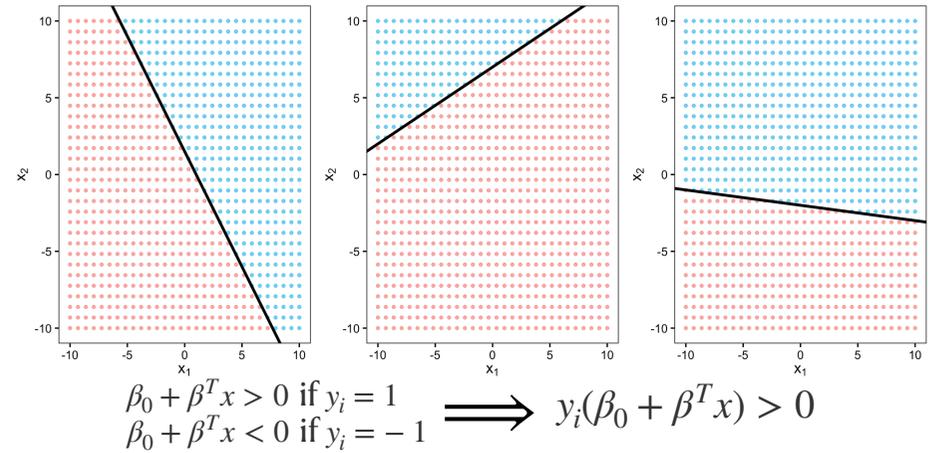
4

Hyperplanes Divide the Space in Half



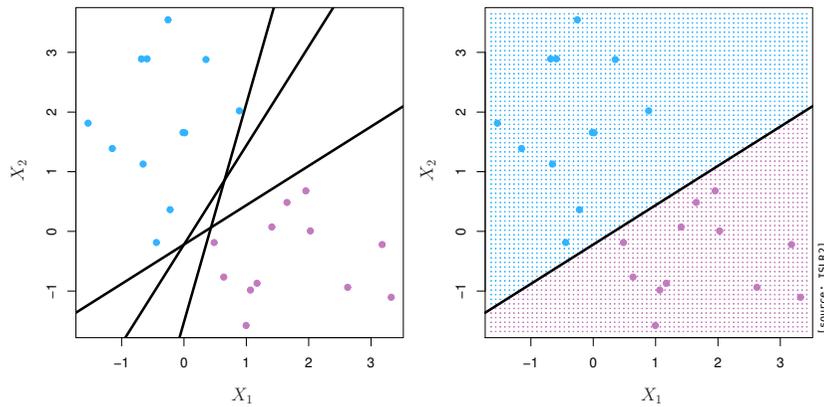
5

Hyperplanes Divide Spaces in Half



6

Hyperplanes Divide Spaces in Half



if a separating hyperplane exists, there will be ∞ many of them
 how do we choose just one?

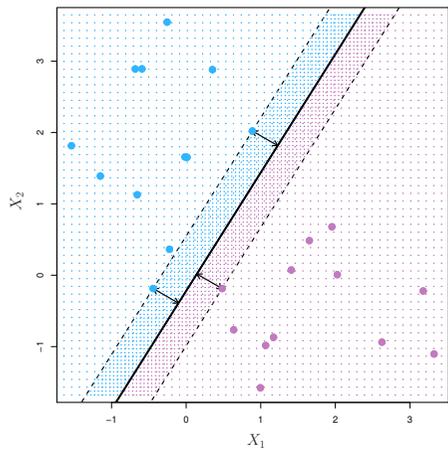
7

Classification



8

Maximal Margin Classifier



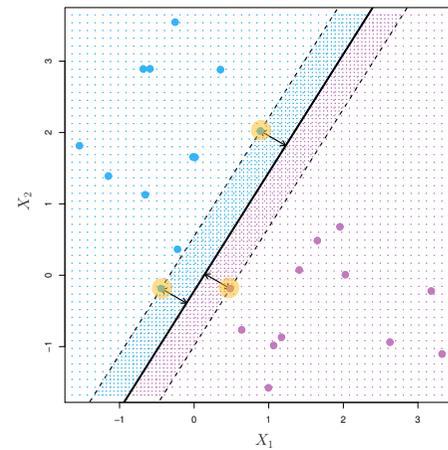
[source: ISLUP2]

choose the separating hyperplane that's furthest from the training examples



9

Maximal Margin Classifier



[source: ISLUP2]

support vectors
the closest points from both classes

the margin
the distance from hyperplane to support vectors

10

Maximal Margin Classifier: The Math

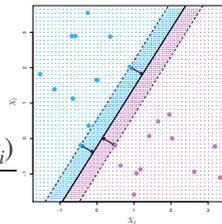
The maximal margin classifier solves a constrained optimization problem:

$$\begin{aligned} & \max_{\beta_0, \beta_1, \dots, \beta_p} M \\ & \text{subject to:} \\ & \|\beta\| = 1 \\ & y_i(\beta_0 + \beta^T x_i) \geq M, \quad \forall i = 1, \dots, n \end{aligned}$$

ensured each observation is on the correct side of the hyperplane and at least a distance M from the hyperplane, i.e., M is the margin of the hyperplane

distance between x_i and line where $y_i(\beta_0 + \beta^T x_i) = M$

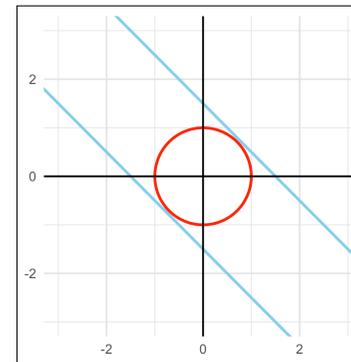
$$\|\beta\| = \sqrt{\sum_{j=1}^p \beta_j^2} \text{ is the Euclidean norm of } \beta$$



11

What is a Constrained Optimization Problem?

Optimize $f(x, y)$ subject to $g(x, y) = k$



$$\begin{aligned} f(x, y) &= 2x + y \\ g(x, y) &= x^2 + y^2 = 1 \end{aligned}$$

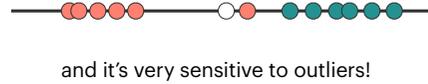
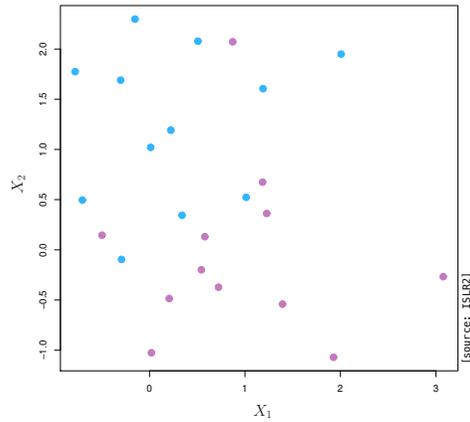
$$\begin{aligned} & \max_{\beta_0, \beta_1, \dots, \beta_p} M \\ & \text{subject to:} \\ & \|\beta\| = 1 \\ & y_i(\beta_0 + \beta^T x_i) \geq M \end{aligned}$$



12

The Non-Separable Case

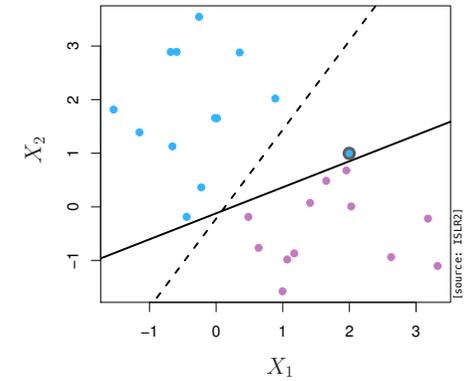
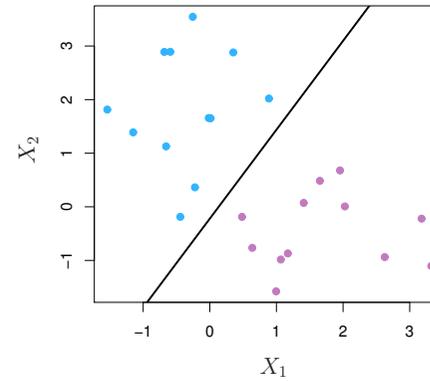
the optimization problem for the maximal margin classifier often has no solution with $M > 0$



13

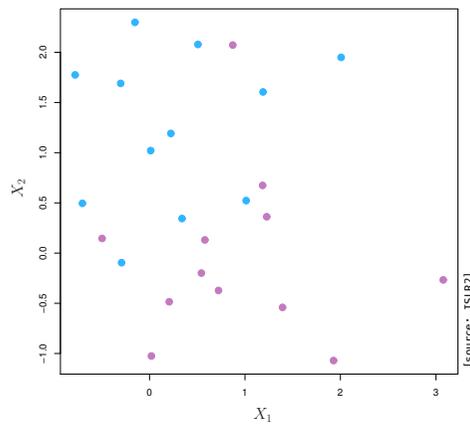
The Non-Separable Case

even if the data are separable, they are sometimes noisy
 \Rightarrow poor solution for the maximum margin classifier



14

The Non-Separable Case

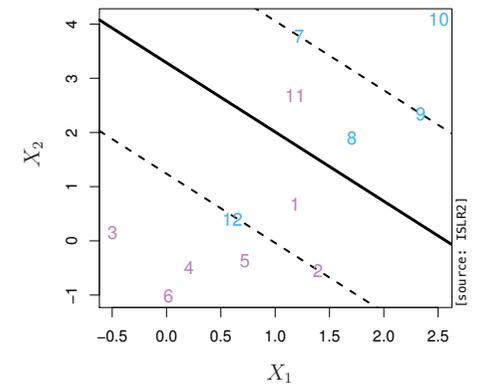
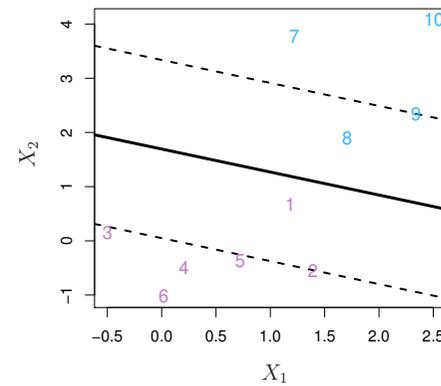


support vector classifier:
 using a so-called **soft margin**
 to **almost** separate the classes

15

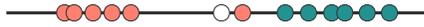
Support Vector Classifier

allows us to classify data that is **not linearly separable**



16

Bias Variance Trade Off



support vector classifier:
how do we choose the soft margin?
→ **cross validation!**

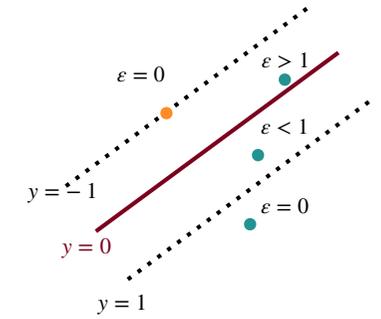
17

Support Vector Classifier

$$\max_{\beta_0, \beta_1, \dots, \beta_p, \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n} M$$

subject to:

$$\begin{aligned} \|\beta\| &= 1 \\ y_i(\beta_0 + \beta^T x_i) &\geq M(1 - \varepsilon_i) \\ \varepsilon_i &\geq 0, \sum_{i=1}^n \varepsilon_i \leq C \end{aligned}$$



$\varepsilon_1, \dots, \varepsilon_n$ are slack variables where $\varepsilon_i = 0$ means i^{th} observation is on correct side of margin
 < 1 means i^{th} observation is on wrong side of margin
 > 1 means i^{th} observation is on wrong side of hyperplane

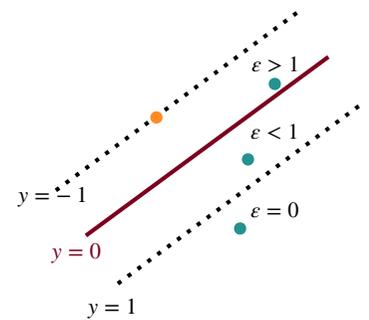
18

Support Vector Classifier

$$\max_{\beta_0, \beta_1, \dots, \beta_p, \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n} M$$

subject to:

$$\begin{aligned} \|\beta\| &= 1 \\ y_i(\beta_0 + \beta^T x_i) &\geq M(1 - \varepsilon_i) \\ \varepsilon_i &\geq 0, \sum_{i=1}^n \varepsilon_i \leq C \end{aligned}$$



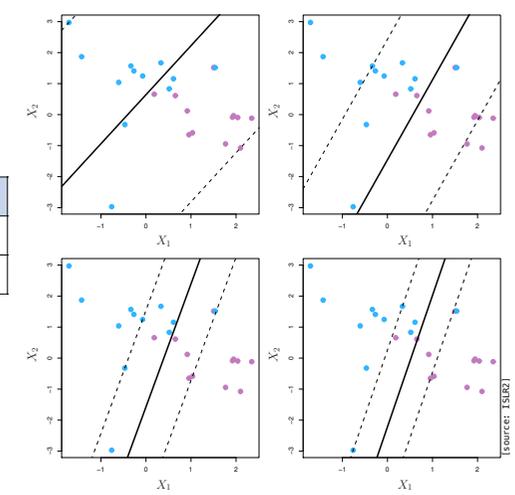
C is the tuning parameter/penalty on error:
 $C = 0$ implies maximal margin hyperplane (superposed it exists)
 $C > 0$ is the total violations to the margin that we can tolerate
 \implies max C observations can be on the wrong side of hyperplane

19

Support Vector Classifier

C : penalty on error

	Regularization	Margins	Bias/Variance
Small C	more	wider	prone to underfitting
Large C	less	narrower	prone to overfitting



20

But what if data looks like this?



21

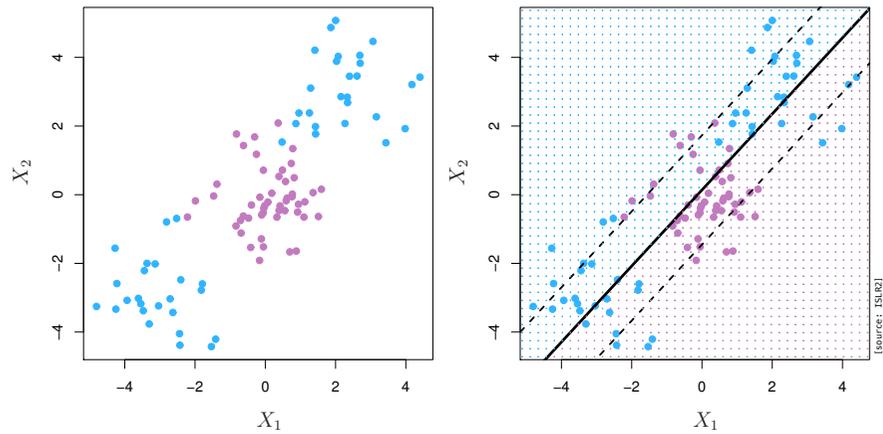
Support Vector Machines



Support Vector Machine (SVM) use Kernel Functions to systematically find Support Vector Classifiers in higher dimensions

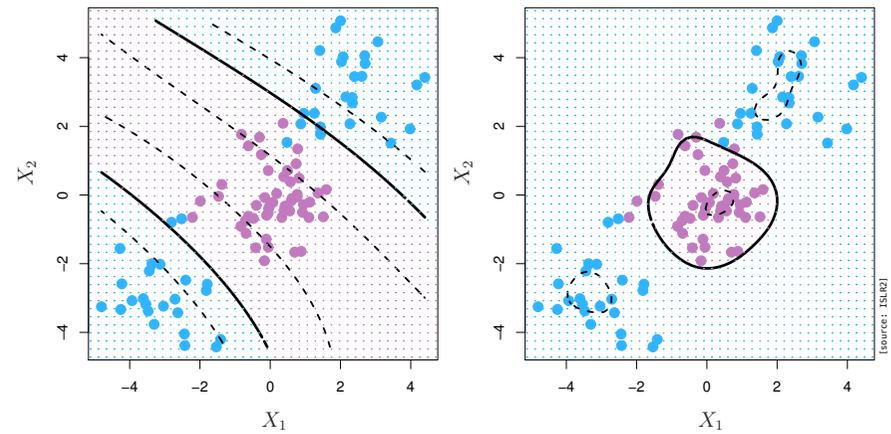
22

Not Linearly Separable Even With Error



23

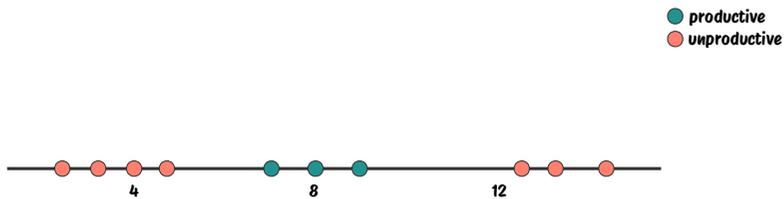
The Kernel Trick



24

The Kernel Trick

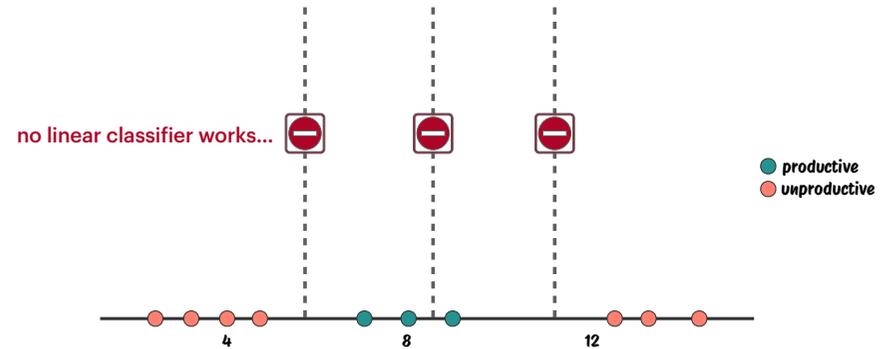
my productivity based on hours of sleep



25

The Kernel Trick

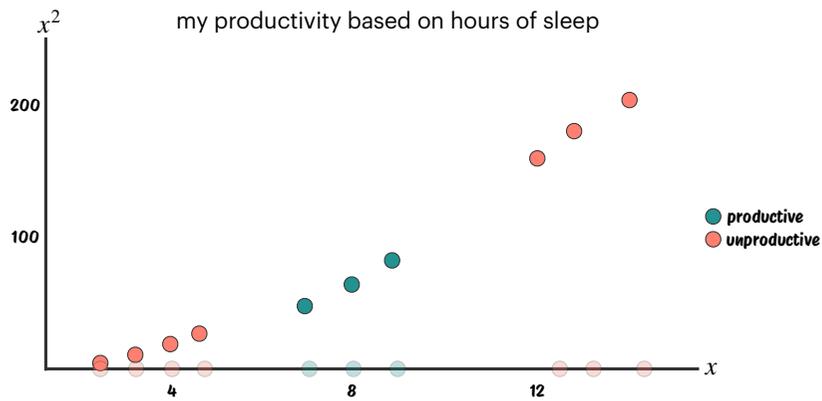
my productivity based on hours of sleep



26

The Kernel Trick

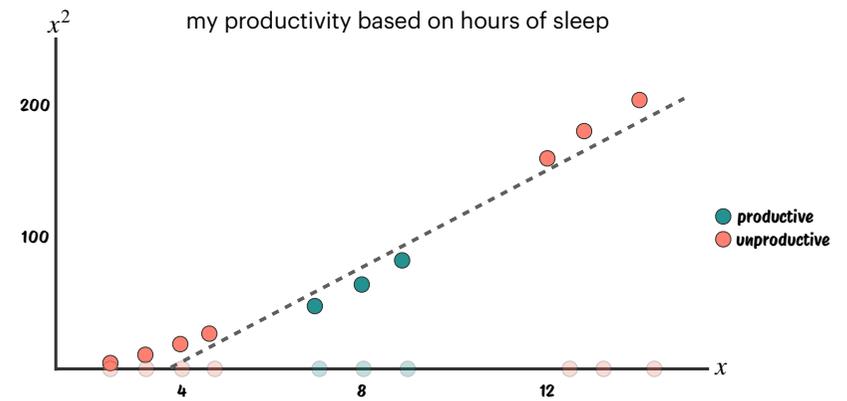
my productivity based on hours of sleep



27

The Kernel Trick

my productivity based on hours of sleep

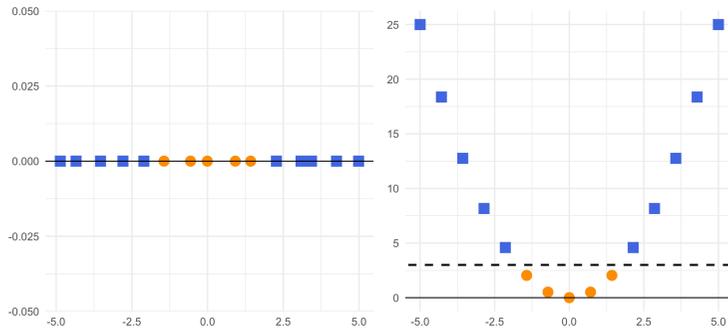


28

The Kernel Trick

what is an SVM Kernel?

A function that computes the relationship between vectors in multiple dimensions (without actually having to calculate the coordinates for those dimensions)



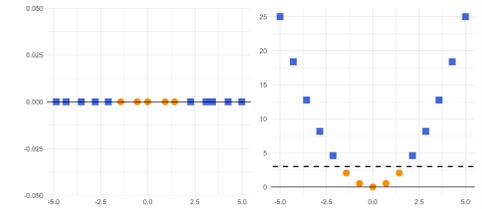
29

The Polynomial Kernel

The **Polynomial Kernel** in the previous sleep vs. productivity example

$$K(a, b) = (a \cdot b + r)^d \quad \text{where } r \text{ is the coefficient and } d \text{ the degree of polynomial (determined by cross validation)}$$

Kernel different observations



30

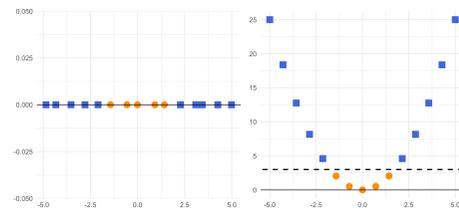
The Polynomial Kernel

The **Polynomial Kernel** in the previous sleep vs. productivity example

$$K(a, b) = (a \cdot b + r)^d \quad \text{where } r \text{ is the coefficients and } d \text{ the degree}$$

we set $r = \frac{1}{2}$ and $d = 2$:

$$\begin{aligned} (a \cdot b + \frac{1}{2})^2 &= (a \cdot b + \frac{1}{2})(a \cdot b + \frac{1}{2}) \\ &+ a^2b^2 + \frac{1}{2}ab + \frac{1}{2}ab + \frac{1}{4} \\ &= \boxed{ab + a^2b^2 + \frac{1}{4}} \end{aligned}$$



31

The Polynomial Kernel

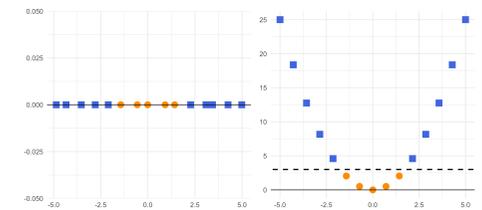
The **Polynomial Kernel** in the previous sleep vs. productivity example

$$K(a, b) = (a \cdot b + r)^d \quad \text{where } r \text{ is the coefficients and } d \text{ the degree}$$

we set $r = \frac{1}{2}$ and $d = 2$:

$$\begin{aligned} (a \cdot b + \frac{1}{2})^2 &= (a \cdot b + \frac{1}{2})(a \cdot b + \frac{1}{2}) \\ &+ a^2b^2 + \frac{1}{2}ab + \frac{1}{2}ab + \frac{1}{4} \\ &= \boxed{ab + a^2b^2 + \frac{1}{4}} = \underbrace{(a, a^2, \frac{1}{2}) \cdot (b, b^2, \frac{1}{2})}_{\text{dot product}} \end{aligned}$$

dot product
gives us the high dimensional coordinates for the data

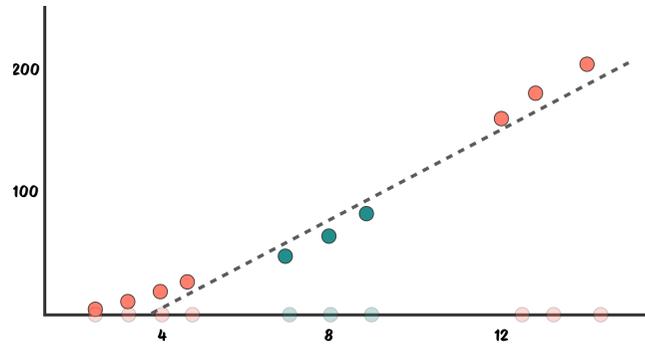


32

The Polynomial Kernel



$$(a, a^2, \frac{1}{2}) \cdot (b, b^2, \frac{1}{2})$$

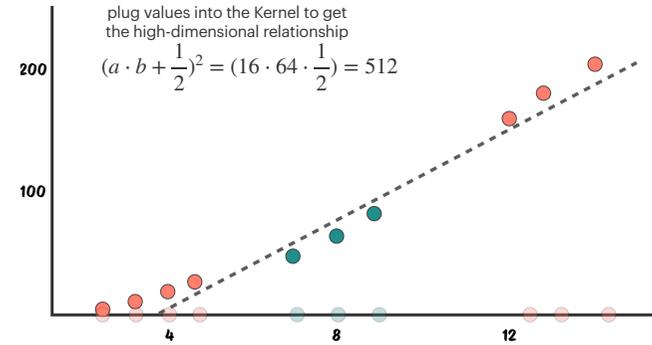


33

The Polynomial Kernel

A function that computes the relationship between vectors in multiple dimensions
(without actually having to calculate the coordinates for those dimensions)

example: $a = 4, b = 8$
 plug values into the Kernel to get
 the high-dimensional relationship
 $(a \cdot b + \frac{1}{2})^2 = (16 \cdot 64 \cdot \frac{1}{2}) = 512$



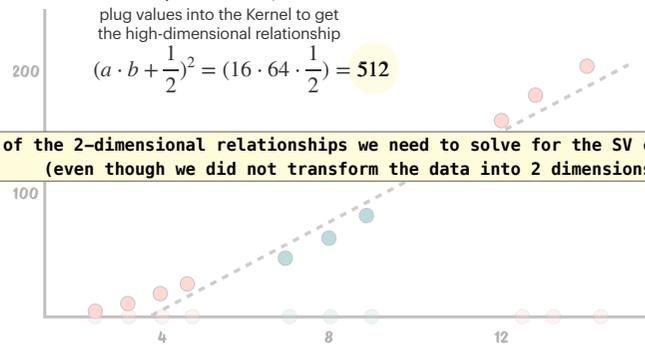
34

The Polynomial Kernel

A function that computes the relationship between vectors in multiple dimensions
(without actually having to calculate the coordinates for those dimensions)

example: $a = 4, b = 8$
 plug values into the Kernel to get
 the high-dimensional relationship
 $(a \cdot b + \frac{1}{2})^2 = (16 \cdot 64 \cdot \frac{1}{2}) = 512$

one of the 2-dimensional relationships we need to solve for the SV classifier
(even though we did not transform the data into 2 dimensions)



35

The Radial Kernel (RBF)

The **Radial Kernel**

$$K(a, b) = e^{-\gamma(a-b)^2}$$

projects to infinite dimensional space
works similar to nearest neighbors classifier

36

The Radial Kernel (RBF)

The **Radial Kernel**

$K(a, b) = e^{-\gamma(a-b)^2}$ projects to **infinite dimensional space**
works similar to nearest neighbors classifier

the amount of influence one observation has on another is a function of the squared distance

$$K(a, b) = e^{-\gamma(a-b)^2}$$

γ scales the squared distance to determine the strength of influence (determined by **cross validation**)



$$K(a, b) = e^{-\gamma(a-b)^2} = \text{high dimensional relationship}$$

37

The Radial Kernel (RBF)

The **Radial Kernel**

$K(a, b) = e^{-\gamma(a-b)^2}$ projects to **infinite dimensional space**
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b + r)^d$$

$$\text{set } r = 0 \implies (a \cdot b)^d = a^d \cdot b^d$$

$$\text{set } d = 1 \implies ab = (a) \cdot (b)$$



38

The Radial Kernel (RBF)

The **Radial Kernel**

$K(a, b) = e^{-\gamma(a-b)^2}$ projects to **infinite dimensional space**
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b + r)^d$$

$$\text{set } r = 0 \implies (a \cdot b)^d = a^d \cdot b^d$$

$$\text{set } d = 1 \implies ab = (a) \cdot (b)$$

$$\text{set } d = 2 \implies a^2b^2 = (a^2) \cdot (b^2)$$



39

The Radial Kernel (RBF)

The **Radial Kernel**

$K(a, b) = e^{-\gamma(a-b)^2}$ projects to **infinite dimensional space**
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b + r)^d$$

$$\text{set } r = 0 \implies (a \cdot b)^d = a^d \cdot b^d$$

$$\text{set } d = 1 \implies ab = (a) \cdot (b)$$

$$\text{set } d = 2 \implies a^2b^2 = (a^2) \cdot (b^2)$$

$$\text{set } d = 3 \implies a^3b^3 = (a^3) \cdot (b^3)$$



40

The Radial Kernel (RBF)

The **Radial Kernel**

$$K(a, b) = e^{-\gamma(a-b)^2}$$

projects to infinite dimensional space
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b + r)^d$$

$$\text{set } r = 0 \implies (a \cdot b)^d = a^d \cdot b^d$$

$$\text{set } d = 1 \implies (a) \cdot (b)$$

$$\text{set } d = 2 \implies (a^2) \cdot (b^2)$$

$$\text{set } d = 3 \implies (a^3) \cdot (b^3)$$



we stay in same dimension when $r = 0$,
but what if we took these polynomials as a sum?

41

The Radial Kernel (RBF)

The **Radial Kernel**

$$K(a, b) = e^{-\gamma(a-b)^2}$$

projects to infinite dimensional space
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b)^d$$

$$ab + a^2b^2 = (a, a^2)(b, b^2)$$



42

The Radial Kernel (RBF)

The **Radial Kernel**

$$K(a, b) = e^{-\gamma(a-b)^2}$$

projects to infinite dimensional space
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b)^d$$

$$ab + a^2b^2 + a^3b^3 = (a, a^2, a^3)(b, b^2, b^3)$$



43

The Radial Kernel (RBF)

The **Radial Kernel**

$$K(a, b) = e^{-\gamma(a-b)^2}$$

projects to infinite dimensional space
works similar to nearest neighbors classifier

we can use the Polynomial Kernel to get the intuition behind how Radial Kernel works in infinite dimensions

$$K(a, b) = (a \cdot b)^d$$

$$ab + a^2b^2 + a^3b^3 + \dots + a^\infty b^\infty = (a, a^2, a^3, \dots, a^\infty)(b, b^2, b^3, \dots, b^\infty)$$

take sum for infinite terms gives dot product with infinite dimensions!

44

The Radial Kernel: Taylor Series Expansion

$$K(a, b) = e^{-\gamma(a-b)^2} = e^{-\gamma(a^2+b^2-2ab)} = e^{-\gamma(a^2+b^2)} e^{\gamma 2ab}$$

$$\text{set } \gamma = \frac{1}{2} \implies e^{-\frac{1}{2}\gamma(a^2+b^2)} e^{ab} \text{ Taylor expansion of this term}$$

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots + \frac{f^{(\infty)}(a)}{\infty!}(x-a)^\infty$$

$$e^x = e^a + \frac{e^a}{1!}(x-a) + \frac{e^a}{2!}(x-a)^2 + \frac{e^a}{3!}(x-a)^3 + \dots + \frac{e^a}{\infty!}(x-a)^\infty, \text{ around } a=0 \text{ we get}$$

$$e^x = 1 + \frac{1}{1!}x + \frac{1}{2!}x^2 + \frac{1}{3!}x^3 + \dots + \frac{1}{\infty!}x^\infty, \text{ so for the point } ab \text{ we get}$$

$$e^{ab} = 1 + (ab) + \frac{(ab)^2}{2!} + \frac{(ab)^3}{3!} + \dots + \frac{(ab)^\infty}{\infty!}$$

each term contains Polynomial Kernel with $r=0$ and d from 0 to $d=\infty$



45

The Radial Kernel: Taylor Series Expansion

$$K(a, b) = e^{-\gamma(a-b)^2} = e^{-\gamma(a^2+b^2-2ab)} = e^{-\gamma(a^2+b^2)} e^{\gamma 2ab}$$

$$\text{set } \gamma = \frac{1}{2} \implies e^{-\frac{1}{2}\gamma(a^2+b^2)} e^{ab} \text{ Taylor expansion of this term}$$

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots + \frac{f^{(\infty)}(a)}{\infty!}(x-a)^\infty$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^\infty}{\infty!}$$

Radial Kernels have coordinates for infinite dimensions!

$$e^{ab} = 1 + (ab) + \frac{1}{2!}(ab)^2 + \frac{1}{3!}(ab)^3 + \dots + \frac{1}{\infty!}(ab)^\infty$$

$$a^0b^0 + a^1b^1 + a^2b^2 + a^3b^3 + \dots + a^\infty b^\infty = (1, a, a^2, a^3, \dots, a^\infty)(1, b, b^2, b^3, \dots, b^\infty)$$



46

The Radial Kernel: Taylor Series Expansion

$$K(a, b) = e^{-\gamma(a-b)^2} = e^{-\gamma(a^2+b^2-2ab)} = e^{-\gamma(a^2+b^2)} e^{\gamma 2ab}$$

$$\text{set } \gamma = \frac{1}{2} \implies e^{-\frac{1}{2}\gamma(a^2+b^2)} e^{ab}$$

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots + \frac{f^{(\infty)}(a)}{\infty!}(x-a)^\infty$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^\infty}{\infty!}$$

Radial Kernels have coordinates for infinite dimensions!

$$e^{ab} = 1 + (ab) + \frac{1}{2!}(ab)^2 + \frac{1}{3!}(ab)^3 + \dots + \frac{1}{\infty!}(ab)^\infty$$

$$e^{ab} = \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \sqrt{\frac{1}{3!}}a^3, \dots, \sqrt{\frac{1}{\infty!}}a^\infty\right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \sqrt{\frac{1}{3!}}b^3, \dots, \sqrt{\frac{1}{\infty!}}b^\infty\right)$$



47

The Radial Kernel: Taylor Series Expansion

$$K(a, b) = e^{-\gamma(a-b)^2} = e^{-\gamma(a^2+b^2-2ab)} = e^{-\gamma(a^2+b^2)} e^{\gamma 2ab}$$

$$\text{set } \gamma = \frac{1}{2} \implies e^{-\frac{1}{2}\gamma(a^2+b^2)} e^{ab}$$

$$e^{-\frac{1}{2}(a-b)^2} = e^{-\frac{1}{2}(a^2+b^2)} \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \dots\right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \dots\right)$$

$$\text{let } s = \sqrt{e^{-\frac{1}{2}(a^2+b^2)}}$$

$$e^{-\frac{1}{2}(a-b)^2} = \left(s, s\sqrt{\frac{1}{1!}}a, s\sqrt{\frac{1}{2!}}a^2, \dots, s\sqrt{\frac{1}{\infty!}}a^\infty\right) \cdot \left(s, s\sqrt{\frac{1}{1!}}b, s\sqrt{\frac{1}{2!}}b^2, \dots, s\sqrt{\frac{1}{\infty!}}b^\infty\right)$$

$$e^{ab} = \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \sqrt{\frac{1}{3!}}a^3, \dots, \sqrt{\frac{1}{\infty!}}a^\infty\right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \sqrt{\frac{1}{3!}}b^3, \dots, \sqrt{\frac{1}{\infty!}}b^\infty\right)$$



48

The Radial Kernel



$$K(a, b) = e^{-\gamma(a-b)^2} = e^{-\gamma(a^2+b^2-2ab)} = e^{-\gamma(a^2+b^2)}e^{\gamma 2ab}$$

$$\text{set } \gamma = \frac{1}{2} \implies e^{-\frac{1}{2}\gamma(a^2+b^2)}e^{ab}$$

$$e^{-\frac{1}{2}(a-b)^2} = e^{-\frac{1}{2}(a^2+b^2)} \left(1, \sqrt{\frac{1}{1!}}a, \sqrt{\frac{1}{2!}}a^2, \dots\right) \cdot \left(1, \sqrt{\frac{1}{1!}}b, \sqrt{\frac{1}{2!}}b^2, \dots\right)$$

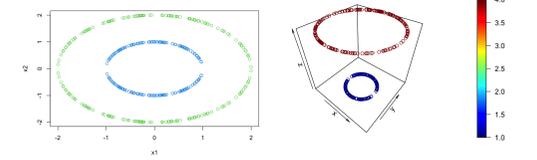
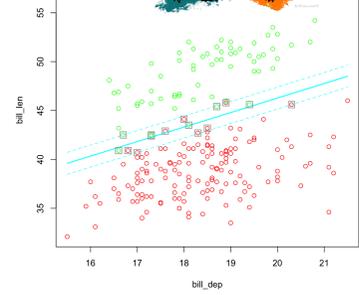
$$\text{let } s = \sqrt{e^{-\frac{1}{2}(a^2+b^2)}}$$

$$e^{-\frac{1}{2}(a-b)^2} = \left(s, s\sqrt{\frac{1}{1!}}a, s\sqrt{\frac{1}{2!}}a^2, \dots, s\sqrt{\frac{1}{\infty!}}a^\infty\right) \cdot \left(s, s\sqrt{\frac{1}{1!}}b, s\sqrt{\frac{1}{2!}}b^2, \dots, s\sqrt{\frac{1}{\infty!}}b^\infty\right)$$

the Radial Kernel is equal to a Dot Product that has coordinates for an infinite number of dimensions.

This Week's Practical

Support Vector Machines



SVM classification plot

