Chapter 6 Kernel method
We considered linear parametric models in which of
the form of mapping y(*,w) is governed by w
of adaptive (trainable) parameters
Training set is used either to find a point estimate
of w or to determine a posterior distribution of w
There is a class of ML in which the training tata
points are kept and used during the prediction phase
E.g. K-NN and Parzen probability density model

They	require	a	metri	c to	be	defined	that	: meo	lsuhes
the	similarity	१ र्ठ	onj	, two	r vec	tors	in in	put s	pace
Linea	r paran	netric	mode	els ca	n be	writte	en in	a	dual
repres	entation.	In	this	form	pred	ictions	are	made	using
linear	combina	ation	రు	kernel	funct	cions ,	which	are	calcu-
lated	using.	the	train	ing do	sta p	oints			
This	works	wel)	when	the	mode	uses	a	fixed	non -
llnear	mappin	g to	o tra	nsform	the	input	tata	into	a

 $K(\mathcal{X}, \mathcal{X}') = \overline{\Phi}(\mathcal{X})^{\mathsf{T}} \overline{\Phi}(\mathcal{X}')$

The simplest kernel function comes from using the identity mapping for the feature transformation $\overline{\Phi}(x) = x$. In this case,

 $K(\mathcal{X},\mathcal{X}') = \mathcal{X}^{\mathsf{T}}\mathcal{X}'$

This called the linear kernel.

Kernel trick (Kernel substitution)

When	an	algor	ithm	uses	inpo	九	vector	rs	only	th	rough	dot
produc	ets,	replo	re	the	40t	pro	oduct	W	ith	a	kernel	
funct	ion	to	make	e it	; Wi	ork	in	a	nonli	near	feat	ture
space												

Dual representation
Consider a linear regression model with regularized
SSE error function given by
$J(w) := \frac{1}{2} \sum_{n=1}^{\infty} \frac{1}{2} w^{T} \Phi(x_{n}) - t_{n} + \frac{1}{2} w^{T} w (6.2)$
where 720 .
Set $\nabla_{w} J(w) = \Theta$. Then we see that the solution for
Wi takes the form of a linear combination of the
vectors $\Phi(X_n)$ with coefficients that are function of w
$W = -\frac{1}{7}\sum_{n=1}^{N} \left\{ w^{T} \overline{\Phi} (X_{n}) - t_{n} \right\} \overline{\Phi} (X_{n}) = \sum_{n=1}^{N} \alpha_{n} \overline{\Phi} (X_{n}) = \overline{\Phi} \alpha_{n} (6.3)$ $M \times N$

where
$$\overline{\Phi}$$
 is the design matrix and $\Delta = (\Delta_1, ..., \Delta_N)^T$ with
 $\Delta_n = -\frac{1}{7} \langle W^T \overline{\Phi} (W^n) - t_n \rangle^T$
We will reformulate the least squares algorithm in
terms of the parameter vector Δ_r which leads to
 Δ dual representation of the algorithm
Substitute $W = \overline{\Phi}^T \Delta_r$ into $J(W)$ (6.2)
 $J(\Delta) = \frac{1}{2} \Delta^T \overline{\Phi} \overline{\Phi}^T \overline{\Phi} \overline{\Phi}^T \Delta_r - \Delta^T \overline{\Phi} \overline{\Phi}^T \overline{\Phi} + \frac{1}{2} \overline{\Phi}^T \overline{\Phi} \overline{\Phi}^T \Delta_r$
where $\underline{\Psi} = (t_1, t_2, ..., t_N)^T$.

Now we define the gram matrix
$$k = \overline{\Phi} \overline{\Phi}^{T}$$
 which is
an $N \times N$ symmetric matrix with elements
 $k_{nm} = \overline{\Phi}(\mathcal{X}_{n})^{T} \overline{\Phi}(\mathcal{X}_{m}) = K(\mathcal{X}_{n}, \mathcal{X}_{m})$
In terms of the gram matrix, SSE error can be written
as
 $J(\alpha) = \frac{1}{2} \alpha^{T} K K \alpha - \alpha^{T} K \# + \frac{1}{2} \#^{T} \# + \frac{\pi}{2} \alpha^{T} K \alpha$
Using (6.3) to eliminate wy from (6.4) and solving.
for α we obtain
 $\alpha = (K + \pi I_{N})^{T} \#$

For linear regression model, we obtain the following.
prediction for a new input
$$\chi$$
 a
 $y(\chi) = W^T \overline{\Phi}(\chi) = a^T \overline{\Phi} \overline{\Phi}(\chi) = [K(\chi)]^T (K+7) I_N]^T \pm (6.9)$
where $|K(\chi)| = W^T \overline{\Phi}(\chi) = K(\chi_n, \chi) = \overline{\Phi}(\chi_n) \overline{\Phi}(\chi)$
for $n = 1, 2, ... N$
The dual formulation allows the solution to the least
square problem to be expressed by $K(\chi, \chi')$
Thus we see that the formulation allows the solution
to the least square problem to be expressed entirely in

terms	र्भ	the	kerr	e	fund	tion	K (%	,*')	(du	ما	formul	ation)
os can	be	writ	ten	as	a	linear		mbinat	tion	q	the	feature
vector	₫৻ж⟩). S	o u	re	Con	recov	ver	the	origi	nal	paro	meter
vector	W.											

Remark

-	То	deter	mine	W,	we	need	OCM3) calcu	alations	
_	In	the	fual	formul	ation,	to	Jetemin	e 🔊	, we ne	et O(N ³)
	The	dual	for	mulatio	on us	es oi	nly the	kerne	el func	tion,
	j.e.	We	40	not	veeq	to c	ompute	the	feature	vector
	Q (M)) dir	ectly							
	We	con	word	k in	high	(even	insinit	e) dim	ensional	space
	with	out	extr	~ co	st (all t	hrough	the	remel)	

기준 basis function 방법
1. Fix a basis function
$$\overline{\Phi}(x) = (\phi_0(x), \phi_1(x), \dots, \phi_{M+1}(x))^T$$

2. Model $Y(x) = W^T \overline{\Phi}(x)$ with weight vector W
3. Finding W minimizing the cost function below
 $J(w) = \frac{1}{2} \sum_{n=1}^{V} \frac{1}{2} W^T \overline{\Phi}(x_n) - t_n \frac{1}{2} + \frac{1}{2} W^T W$
where $71 \ge 0$.

kernel	trick (dual representation)
1. Fix	a kernel function KC%, %')
2. Model	$Y(x) = k(x)^{T} \otimes where k(x) = (k_1(x), k_2(x),, k_N(x))^{T}$
and kic	$(\mathscr{X}) = \mathcal{K}(\mathcal{X}_{1}, \mathcal{X})$
3. Finding	as minimizing the cost function below
20	ころ) = ビューズ KKの - の KH + ビモ + ビタン Kの
where 1	$c_{nm} = K(\mathcal{H}_{n}, \mathcal{H}_{m}), \ \mathcal{J} \geq 0$

6.2	Constructing	kernels				
In oi	rder to use	kernel s	ubstitutio	on, we	neet to	construct
valid	kernel fu	nctions.				
One	approach is	to ch	oose a	feature	space	mapping
₫(*)	and then	use this	to fin	d the	correspond	ing kernel
	× ۲.	$x') := \overline{\mathbb{Q}}(x)$	$\Phi(x) = $	$\sum_{i=1}^{M} \phi_i(x) \varphi_i$	⁵ . (*')	
where	$\phi_{x}(x)$ are	basis func	tions.			



Figure 6.1 Illustration of the construction of kernel functions starting from a corresponding set of basis functions. In each column the lower plot shows the kernel function k(x, x') defined by (6.10) plotted as a function of x, where x' is given by the red cross (×), while the upper plot shows the corresponding basis functions given by polynomials (left column), 'Gaussians' (centre column), and logistic sigmoids (right column).

Alternative	e approc	ich is	to	built 1	cernel f	unctio	n dir	ectly
without	thinking	about	the	feature	e mapp	ing.	€ (%)	
flowever,	to be	a valid	kerr	el, the	e funct	ion n	nust	sastis-
fy certain	n conditio	ons: It	shou	19 con	respond	to	a	dot
product	in some	e featu	re sp	pace (even in	finite	dime	nsional)
Mathematic	cally, thi	s means	the	kernel	matrix	k	where	¢
$k_{nm} = k($	(*******)	must	be					
– symm	etric Kn	m = Kmn						

- Positive semi-definite

E.g. Consider the kernel function given by

$$K(\mathcal{X}, \mathcal{Z}) := (\mathcal{X}^{T} \mathcal{Z})^{2}$$
In case of $D = 2$, we can expand out the terms and
identify the corresponding nonlinear feature mapping

$$K(\mathcal{X}, \mathcal{Z}) = (\mathcal{X}^{T} \mathcal{Z})^{2} = (\mathcal{X}_{1} \mathcal{Z}_{1} + \mathcal{X}_{2} \mathcal{Z}_{2})^{2}$$

$$= \chi_{1}^{2} \mathcal{Z}_{1}^{2} + 2\chi_{1} \mathcal{Z}_{1} \chi_{2} \mathcal{Z}_{2} + \chi_{2}^{2} \mathcal{Z}_{2}^{2}$$

$$= (\chi_{1}^{2}, J_{2} \chi_{1} \chi_{2}, \chi_{2}^{2}) (\mathcal{Z}_{1}^{2}, J_{2} \mathcal{Z}_{2}, \mathcal{Z}_{2}^{2})^{T}$$

$$= \overline{Q}(\mathcal{X})^{T} \overline{Q}(\mathcal{Z})$$
Feature mapping $\overline{Q}(\mathcal{X}) = (\chi_{1}^{2}, J_{2} \chi_{1} \chi_{2}, \chi_{2}^{2})^{T}$

Here are	the main	properties	that	allow	combining	simpler
kernels to	build mor	e poverful	one:			
	Techniques for	Constructing New Kern	els.			
	Given valid kerne be valid:	els $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$), the followin	ng new kernels	will also	
		$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x})$	$\mathbf{x}')$		(6.13)	
		$k(\mathbf{x},\mathbf{x}') = f(\mathbf{x})k_1$	$(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$		(6.14)	
		$k(\mathbf{x},\mathbf{x}') = q(k_1(\mathbf{x}))$	$(\mathbf{x}'))$		(6.15)	
		$k(\mathbf{x},\mathbf{x}') = \exp\left(k_1\right)$	$(\mathbf{x},\mathbf{x}'))$		(6.16)	
		$k(\mathbf{x},\mathbf{x}') ~=~ k_1(\mathbf{x},\mathbf{x})$	$k')+k_2(\mathbf{x},\mathbf{x}')$)	(6.17)	
		$k(\mathbf{x},\mathbf{x}') ~=~ k_1(\mathbf{x},\mathbf{x})$	$k')k_2({f x},{f x}')$		(6.18)	
		$k(\mathbf{x},\mathbf{x}') = k_3 \left(\boldsymbol{\phi}(\mathbf{x},\mathbf{x}') \right)$	$\mathbf{x}), oldsymbol{\phi}(\mathbf{x}'))$		(6.19)	
		$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^{\mathrm{T}} \mathbf{A} \mathbf{x}'$			(6.20)	
		$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a,$	$\mathbf{x}_a') + k_b(\mathbf{x}_b,$	$\mathbf{x}_b')$	(6.21)	
		$k(\mathbf{x},\mathbf{x}') \;\;=\;\; k_a(\mathbf{x}_a,$	$\mathbf{x}_a')k_b(\mathbf{x}_b,\mathbf{x}_b')$)	(6.22)	
	where $c > 0$ is a cative coefficients \mathbb{R}^M , A is a symmotry disjoint over their respective.	constant, $f(\cdot)$ is any funct, , $\phi(\mathbf{x})$ is a function from hetric positive semidefinit nt) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, an ive spaces.	ion, $q(\cdot)$ is a p x to \mathbb{R}^M , k_3 te matrix, \mathbf{x}_a a d k_a and k_b and	olynomial with (\cdot, \cdot) is a valid and \mathbf{x}_b are variance re valid kernel f	kernel in bles (not functions	

Polynomial kernel (example)

$$k(x, x') = (x^T x')^2$$

includes only terms of degree 2.

If we add a constant before squaring like

$$K(x',x') = (x''x' + c)^{2} \quad \text{where } c > 0$$
then the feature space becomes richer.
This kernel includes a constant (bias), linear terms (x_{1}, x_{2})
and degree 2 terms $(x_{1}x_{2} \text{ or } x_{1}^{2})$

$$k(x', x') = (x_{1}x'_{1} + x_{2}x'_{2} + c)^{2} = (x_{1}x'_{1})^{2} + 2x_{1}x'_{1}x_{2}x'_{2} + (x_{2}x'_{2})^{2}$$

$$+ 2c x_{1}x'_{1} + 2c x_{2}x'_{2} + c^{2}$$

$$= \overline{\varrho}(x)^{T} \overline{\varrho}(x')$$

$$\overline{\varrho}(x) = (\sqrt{2c} x_{1}, \sqrt{2c} x_{2}, x_{1}^{2}, \sqrt{2c} x_{2}, x_{2}^{2}, c)^{T}$$

M k(水,水') = (メ^T メ')

gives	a	feature	space	with	all	monomials	ир	to
degree	Μ							

One	Co	ommon	ly u	sed	ke	ernel	ìs	the	e (Taussi	an	keme	el	
ca.k.	a	RB	SF ke	ernel).									
			k	¥,⊁′)	=	exp	(<u>(% –</u> 20	- */11 ²)				
This	k	erne)	measi	ires	ha	ow	clo	se	two	vec	tors	×	and	*'
ore.														
k (*, ;	*')	is	close	de de	١	if	¥	and	*	are	sim	ilar		
k (*, ;	∦′১	is	close	क	0	if	¥	and	₩′	are	far	ap	art	
Althou	ıqh	it	look	s li	ke	م	Grau	issian	Pr	robabil	ity	densit	×Y,	here
it i	S	not	used	as	C	r	proba	bility	•					

RBF kernel is a valid kernel $\| x - x' \|^{2} = x^{T} + (x')^{T} x' - 2 x^{T} x'$ So $\mathsf{K}(\mathscr{K},\mathscr{K}') = \exp\left(-\frac{\mathscr{K}^{\mathsf{T}}\mathscr{K}}{2\mathfrak{S}^{2}}\right) \exp\left(-\frac{\mathscr{K}^{\mathsf{T}}\mathscr{K}'}{2\mathfrak{S}^{2}}\right) \exp\left(-\frac{(\mathscr{K}')^{\mathsf{T}}\mathscr{K}'}{2\mathfrak{S}^{2}}\right)$ Beause of (6.14) and (6.16), together with the validity of the linear kernel $k(x, x') = x^T x'$

Remark

- The Graussian kernel corresponds to an infinite - dimensional
feature mapping $\Phi(x)$ (exercise 6.11)
- The Goussian kernel is not restricted to the use of
Euclidean distance. I.e. *** can be replaced with
any nonlinear kernel K(*,*'),
$K(X,X') = exp(-\frac{1}{25^2}(\widehat{K}(X,X) + \widehat{K}(X',X') - 2\widehat{K}(X',X')))$

One	useful	way	to k	wild ke	ernel	is	by usin	g a		
probab	oilistic	genera	tive r	nodel						
Suppos	e we	have	a ge	nerative	, mc	del	PC*).	Then	we	
con	define	ak	emel	05						
			KCX,	*′)= f) (*) P	(*′)				
This	is a	valid	because	: it	CONN	be	written	20	an	inner
produc	£ :									
		KCX,	*′)= ∮		¢) w	here	₫(*)=	PC*)		

Note	e that	it	is	a	very	simple	, featu	re map	ping.	where
each	i input	Ж	is	just	mar	ped	to a	scalar	value	PC*).
In	practice	, mo	re	ad var	ncet	versio	ns use	likelih	00 4 f	unction
or	posteri	or di	strib	utions	0\$	feat	ures			
Γŧ	says	that	X	and	*′	are	similar	· i s	they	both
have	high	pro	babi	ities						

Example: kernel from a Graussian Mixture Model
Assume
$$p(x) = \sum_{k=1}^{k} \pi_k N(x|\mu_k, \Sigma_k)$$

 $p(z=k) P(x|z=k)$
and estimates π_k , μ_k , Σ_k for $k=1,2,...,k$
For each x , calculate the posterior responsibility
 $\phi_k(x) = p(z=k|x) = \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^{k} \pi_j N(x|\mu_j, \Sigma_j)}$
Define $\overline{\Psi}(x) := (\phi_1(x), \phi_2(x), ..., \phi_k(x))^T$ and
 $k(x, x') = \overline{\Psi}(x)^T \overline{\Psi}(x')$

This	kernel	measo	ures the	similarit	y between	two	vectors
based	on	how	responsib	ole each	Craussian	compone	ent is
for	genera	ating	them.			19	

kemels	from mixtur	e models	with lat	ent vari	able
We can	build more	flexible	kernels	by sur	ming over
multiple	components	in a pr	obabilistic	mixture	model
1. Disc	rete case				
Suppose	we have				
- lat	ent variable	ie 11,2,	k }		
- prok	. distribution	for each	compone	ent pc	×Ιλ)
- pric	or over comp	ponents p	>(ふ)		
	The kernel	k(*,*')=	- Σ P(*1) ===1	ριχίλγριλ)

_	This	kerne	is	large	whe	n X	and	*	both	have	high
	likelih	000 (under	the	sar	ve	сомро	nents			
_	Γt	reflec	ts	how	simila	arly	the	mode	l ex	plains	the
	two	inputs									
~	The	laten	t va	ariable	λ	can	be	seen	as	a	latent
	variab	le (e	g.,	cluster	ind	2X	in a	mixt	ure	model)

2.	Conti	nyous	. ca	se							
Ŀf	the	. 10	ctent	V	ariable	Æ	is	conti	nuous	we	replace
the	SUM	\sim	ith	01	integr	al					
		۲C#,#	:') =	SP	(∦\₹)	PC*'1	₹) P((3) 43	£		
whe	re	7	is	a	contin	uous	lat	ent	variable	e.	

An	other	example	of	a	kernel	function	is	the	sigmoidal
kei	mel								
			× ۲.× , »	<') =	tanh (o	. * ^T *′ + b)			
Re	mark								
1	This	function	looks		ke the	activatio)N	function	n in
	neura	l netwo	-k						
~	This	kernel	is n	ot	always	valid			
	Cits	gram	matrix	5	not	guarantee	9	to be	positive
	semi	Jefinite	for	all	a, b	and do	ta	set).	



SVM	W/*+b	
min www.		quadratic objective
st. t; (w ^T *; + b) 2	×=1, N	linear constraints
This is a quadratic optim	ization problem.	
The one method that maxim	nizes the distance	to the
closest data points from	both classes	
- Maximal margin		






$$\frac{\max_{W,b} + cw, b}{W,b} = \max_{W,b} \min_{X \in D} \frac{|w^T + b|}{||w||}$$
s.t.
$$t_{\lambda} (w^T + b) \ge 0 \quad \forall_{\lambda} = 1, 2, ... N$$

Quite complicated. Need to simplify the problem





We need to simplify the constraints





SVM





Which one is better?



SVM

Slack variable
$$\frac{1}{2}$$

- Data points are allowed to be on 'urong side'
of the margin boundary but with a penalty that
increases with the distance from that boundary
 $\frac{1}{2} \| w \|^{2} + C \cdot \sum_{\lambda=1}^{N} \frac{1}{2} \| w \|^{2} + C \cdot \sum_{\lambda=1}^{N} \frac$

$$\begin{array}{rcl} \vdots = 0 & : & connectly & classified \\ 0 < \vdots \leq 1 & : & inside & the margin \\ & & but & on & the & connect & side \\ \vdots > 1 & & y = 0 \\ & & & y = 0 \\ \xi > 1 & & y = 1 \\ \hline \xi > 1 & & y = 1 \\ \xi < 1 & & \\ \xi < 1 & & \\ \xi = 0 \\ \hline \xi = 0 \end{array}$$

$$\begin{array}{c} \min_{\substack{i=1\\ W_i, i \\ W_i$$



Overlapping classes : Hinge loss

$$\begin{array}{c}
\text{min} \quad \frac{1}{2} \parallel w \parallel^{2} + c \cdot \sum_{\lambda=1}^{N} \quad & (\chi^{T} \chi') \rightarrow \quad & (\chi, \chi') \\
\text{st.} \quad t_{\lambda} (w^{T} \chi_{\lambda} + b) \geq 1 - \frac{1}{2} \\
& & \chi_{\lambda} = 1 \dots N \\$$

$$\Rightarrow \qquad \underset{w, s}{\min} \quad \frac{1}{2} \parallel w_{\parallel} \parallel^{2} + C \cdot \sum_{\lambda=1}^{N} s_{\lambda}$$

$$\Rightarrow \qquad \underset{w, s}{\sup} \quad \frac{1}{2} \parallel w_{\parallel} \parallel^{2} + C \cdot \sum_{\lambda=1}^{N} s_{\lambda}$$

$$s.t. \quad s_{\lambda} = \max(0, 1 - t_{\lambda} \gamma(s_{\lambda})) \quad \forall_{\lambda=1}, \dots N$$

where $\gamma(s_{\lambda}) = w^{T} s_{\lambda} + b$

$$= \sum_{\substack{m \in I \\ m \in I}} \min_{j \in I} \frac{1}{j} \lim_{\substack{m \in I \\ n \neq i}} \frac{1}{j} \lim_$$

or min
$$\sum_{k=1}^{N} \max(0, 1-t_{k}\gamma(x_{k})) + \frac{1}{2c} \|W\|^{2}$$

where
$$Y(*_{\lambda}) = W^{T}*_{\lambda} + 1$$

SVM = Linear classifier with Hinge loss

$$\begin{array}{rcl}
& \text{min } \sum_{\lambda=1}^{N} \max\left(0, 1-t_{\lambda} \gamma(x_{\lambda})\right) + \frac{1}{2c} \|W_{\parallel}\|^{2} \\
& \text{where } \gamma(x_{\lambda}) = W^{T} x_{\lambda} + b
\end{array}$$

Kernel SVM





Kernel method 방법 (key idea) - If an algorithm only uses dot products of vectors, then we can replace the dot product with a kernel function ¥^T* K(X, X')Ð $\Phi(x) \Phi(x')$ 유사도 또는 상관 정도 질것 How to construct a valid kernel function?

Greneralized linear regression Linear model: $Y(*) = W^T * + b$ Basis function method 1. Fix a basis function $\Phi(*) = (\phi_0(*), \phi_1(*), \dots, \phi_{M-1}(*))'$ 2. Model $Y(X) = W^T \Phi(X)$ with weight vector W3. Finding W/ minimizing the cost function below $J(w) = \frac{1}{2} \sum_{n=1}^{\infty} \frac{1}{2} w^{T} \Phi(x_{n}) - t_{n} \sum_{n=1}^{\infty} \frac{1}{2} w^{T} w$ where 720. 4. Inference : $y(\hat{x}) = w^T \Phi(\hat{x})$

$$J(w) = \frac{1}{2}(\Phi w - t)^{T}(\Phi w - t) + \frac{2}{2}w^{T}w$$

Set
$$\nabla_{w} J (w) = \emptyset$$
. Then we see

$$w = \sum_{n=1}^{N} \alpha_n \overline{\Phi}(x_n) = \overline{\Phi}^T \alpha_n^{N \times 1} (6.3)$$
where $\overline{\Phi}$ is the design matrix and $(6.3) = (\alpha_1, \dots, \alpha_N)^T$ with
 $\alpha_n = -\frac{1}{2} \langle w^T \overline{\Phi}(x_n) - t_n \rangle$

Substitute
$$w_{l} = \overline{\Phi}^{T} a$$
 into $J(w)$ (6.2)

$$J(a) = \frac{1}{2} a^{T} \overline{\Phi} \overline{\Phi}^{T} \overline{\Phi} \overline{\Phi}^{T} \overline{\Phi} - a^{T} \overline{\Phi} \overline{\Phi}^{T} \overline{\Phi} + \frac{1}{2} \overline{\Phi}^{T} \overline{\Phi} + \frac{1}{2} \overline{\Phi}^{T} \overline{\Phi} \overline{\Phi}^{T} a$$
where $\underline{\Psi} = (\underline{t}_{1}, \underline{t}_{2}, \dots \underline{t}_{N})^{T}$.

$$\overline{\Phi} \overline{\Phi}^{T} = \begin{pmatrix} -\overline{\Phi}(x)^{T} - \\ -\overline{\Phi}(x)^{T} - \\ \vdots \\ -\overline{\Phi}(x)^{T} - \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ \overline{\Phi}(x) & \overline{\Phi}(x) \end{pmatrix} \qquad N \times N$$

kernel trick (dual representation)	
1. Fix a kernel function KC×,×')	
2. Model $Y(x) = k(x)^T \infty$ where $ k(x) = (k_1(x), k_2(x),, k_N(x))$	ſ
and $k_1(x) = k(x_1, x)$	
3. Finding a minimizing the cost function below	
Jca) = farkka - arkt + ft + + farka	
where $K_{nm} = K(x_n, x_m), \ 7 \ge 0$	
4. Inference $\gamma(\hat{x}) = k(\hat{x})^T a = \sum_{i=1}^{N} a_i k(x_i, \hat{x})$	

kemel support
 vector
 machine

 Primial
 problem

 min

$$\frac{1}{2} \parallel w \parallel^2$$
 st.
 $t_{\bar{\lambda}} (w^T *_{\bar{\lambda}} + b) \ge 1$
 $\forall_{\bar{\lambda}} = 1, ... N$

 w.b
 $\frac{1}{2} \parallel w \parallel^2$
 st.
 $t_{\bar{\lambda}} (w^T *_{\bar{\lambda}} + b) \ge 1$
 $\forall_{\bar{\lambda}} = 1, ... N$



Primial problem (hand margin)
min
$$\frac{1}{2} || w ||^2$$
 st. $t_{\lambda} (w^T *_{\lambda} + b) \ge |^{-V_{\lambda}} = 1... N$
Construct Lagrangian
To handle the inequality, introduce Lagrangian multipliers
 $\chi_{\lambda} \ge 0$ and form the Lagrangian:
 $\int (w_{\lambda}, b, \alpha) = \frac{1}{2} || w ||^2 - \sum_{\lambda=1}^{N} \chi_{\lambda} (t_{\lambda} (w^T *_{\lambda} + b) - 1)$
 $\leq \frac{1}{2} || w ||^2$
 $\chi_{\lambda} := (\chi_{\lambda}, ..., \chi_{N})^{T}$

Derive the dual
To obtain the dual, we minimize the Lagrangian w.r.t.
W and b, and maximize w.r.t. ox. I.e.
min max
$$\downarrow (w, b, ox)$$

 $w_{i,b} \neq 20$
 $-\frac{\partial \bot}{\partial w} = \circledast \implies w_{i} = \sum_{\lambda=1}^{N} \alpha_{\lambda} t_{\lambda} \times \sum_{\lambda=1}^{N} -\frac{\partial \bot}{\partial b} = 0 \implies \sum_{\lambda=1}^{N} \alpha_{\lambda} t_{\lambda} = 0$

Substitute these back into the Lagrangian

Eliminating w and b from
$$L(w, b, \sigma)$$
 using these
conditions. Then the gives the dual representation
$$\widehat{J}(\sigma x) := \sum_{\lambda=1}^{N} \alpha_{\lambda} - \frac{1}{2} \sum_{\lambda,j} \alpha_{\lambda} \alpha_{j} t_{\lambda} t_{j} \overset{T}{}_{\lambda} \overset{T}{}_{j}$$

s.t. $\alpha_{\lambda} \ge 0$ $\lambda = 1 \dots N$, $\sum_{\lambda=1}^{N} \alpha_{\lambda} t_{\lambda} = 0$

Dual optimization problem $\max \sum_{\lambda=1}^{N} \alpha_{\lambda} - \frac{1}{2} \sum_{\lambda,j} \alpha_{\lambda} \alpha_{j} t_{\lambda} t_{j} * *_{\lambda} *_{j}$ s.t. $\alpha_{\lambda} \ge 0$ $\lambda = 1 \dots N$, $\sum_{\lambda=1}^{N} \alpha_{\lambda} t_{\lambda} = 0$

Kernel	trick	[Kemel	SVM)			
If th	e tata	is not	linearly	separable	, we	map the
input	data	to a h	igher - dia	nensional	space	using a
basis	function	重(·)	or we	replace	the	tot product
₩ ^T % ′	with a	. kernel	KC#,#). I.e.		
		N S N -	<u> </u>	x. +. +. K(* * *	
	m o S		ZZŃ	J LA J III	· X , ···j /	
	s.t.	x;20	λ=1,N	こ え え え え え え え え え え え え え え え え え え え	$t_{\lambda} = 0$	



SVM: Overlapping classes (reall)

$$\begin{array}{c}
\text{min} \quad \frac{1}{2} \parallel w \parallel^{2} + c \cdot \sum_{\lambda=1}^{N} \frac{1}{\lambda} \quad \text{primal problem} \\
\text{w.s} \quad \frac{1}{2} \parallel w \parallel^{2} + c \cdot \sum_{\lambda=1}^{N} \frac{1}{\lambda} \quad \text{primal problem} \\
\text{st.} \quad t_{\lambda} (w^{T} \times_{\lambda} + b) \geq 1 - \frac{1}{\lambda} \quad \frac{1}{\lambda - 1} \\
\frac{1}{\lambda} \geq 0 \quad \frac{1}{\lambda - 1} \dots \quad N
\end{array}$$

Here C>O controls the trade-off between the slack Variable penalty and the margin

To obtain the dual, we minimize the Lagrangian w.r.t.
W and b, and maximize w.r.t. ok. and us
min max
w.t. ex.,
$$M$$
 and M
where $\mathcal{L}(W, b, t, ex., M)$
 $= \frac{1}{2} ||W||^2 + C \sum_{\lambda=1}^{N} t_{\lambda} (t_{\lambda} Y(x_{\lambda}) - 1 + t_{\lambda}) - \sum_{\lambda=1}^{N} \mathcal{M}_{\lambda} t_{\lambda};$
 $|X_{\lambda} \ge 0 \}$ and $\int \mathcal{M}_{\lambda} \ge 0 \}$ are Lagrange multipliers.

We optimize out w, b and isit





Using these results to eliminate w, b and $1\frac{1}{2}$ from the Lagrangian, we obtain the dual Lagrangian $\widehat{\mathcal{L}}(x) = \sum_{k=1}^{N} \alpha_k - \frac{1}{2} \sum_{\lambda j} \alpha_{\lambda} \alpha_j t_{\lambda} t_{j} *_{\lambda}^{T} *_{j}$

We note that	X,ZO	and x	$\lambda = C - \lambda$	۸ _λ	
Since Mi 20,	, we ha	ve K _x E	С		
Therefore we	have to	minimize	J COR)	w.r.t	1 Kit
subject to					
	04	≪ [×] ∠ C	for $\lambda = 1, N$		
	$\sum_{n=1}^{N} \alpha_{n}$	$k_{\lambda} t_{\lambda} = 0$			
	۱= ۸				

Kernel SVM: Overlapping classes

$$\max_{\substack{N \\ N \\ N}} \sum_{\lambda=1}^{N} \alpha_{\lambda} - \frac{1}{2} \sum_{\lambda,j} \alpha_{\lambda} \alpha_{j} t_{\lambda} t_{j} K(\mathcal{K}_{\lambda}, \mathcal{K}_{j})$$
s.t. $0 \leq \kappa_{\lambda} \leq C$

$$\forall_{\lambda} = 1, \dots, \qquad \sum_{\lambda=1}^{N} \alpha_{\lambda} t_{\lambda} = 0$$