



PARENTS: TALK TO YOUR
KIDS ABOUT LINUX...
BEFORE SOMEBODY ELSE DOES.

Announcements

From the course reps – Taylor Qin and
Yuchen Fang also see piazza @229

Education info session: FAQ on admin
your degree

After the break:

Week 7 Monday - Easter; Week 8 Monday - ANZAC day: no lecture or tutorial, join another online tutorial the same week.

Topic: How to Admin Your Degree

Date: Friday 1 April 10am

Venue: Rm 1.36 Birch Building (35) or

<https://anu.zoom.us/j/88641270633?pwd=ZU9paE1nUzJ0SU90c1krRnBJMWpGUT09>

Facilitator: Paul Dowden, Manager, Student Services, Employability & Experience

About this event

There are a number of important elements CECS coursework students need to be aware of when managing their degree. The session will include:

1. Terminology (program/degree, major/minor; maximum/minimum etc.)
2. Building a degree:
 - Programs and Courses (P&C)
 - units of credit
 - Compulsory courses
 - Program List electives
 - 'Free' electives
 - Minors - Majors/Specialisations what are they and how and when to register
3. Program planning using P&C
4. What to do if it goes wrong!
 - Repeating a course
 - Reduced Study load
 - Program transfers incl. Credit & Exemption

Kernels

Basis functions recap

Non-parametric methods (Chap 2.5)

Dual representation (of linear models, Chap 6.1)

Constructing kernels (Chap 6.2)

Next time: maximum-margin classifiers (Chap 7.1)

Kernel

From Wikipedia, the free encyclopedia

Not to be confused with [Colonel](#).

Kernel may refer to:

Computing [edit]

- Kernel (operating system), the central component of most operating systems
- [Kernel \(image processing\)](#), a matrix used for image convolution
- Compute kernel, in GPGPU programming
- [Kernel method](#), in machine learning
 - In numerical analysis, a subroutine that performs a common numerical operation
 - In particular, a routine that is executed in a vectorized loop; see, e.g., [General-purpose computing on graphics processing units § Kernels](#)
 - Kernelization, a technique for designing efficient algorithms

Look up *kernel* in
Wiktionary, the free
dictionary.

Contents [hide]

- 1 Computing
- 2 Mathematics
 - 2.1 Objects
 - 2.2 Functions
- 3 Natural sciences
- 4 Other uses
- 5 People
- 6 See also

Mathematics [edit]

Objects [edit]

- Kernel (algebra), a general concept that includes:
 - Kernel (linear algebra) or null space, a set of vectors mapped to the zero vector
 - Kernel (category theory), a generalization of the kernel of a homomorphism
 - Kernel (set theory), an equivalence relation: partition by image under a function
 - Difference kernel, a binary equalizer: the kernel of the difference of two functions

Functions [edit]

- Kernel (geometry), the set of points within a polygon from which the whole polygon boundary is visible
- [Kernel \(statistics\)](#), a weighting function used in [kernel density estimation](#) to estimate the probability density function of a random variable
- Integral kernel or [kernel function](#), a function of two variables that defines an integral transform
- Heat kernel, the fundamental solution to the heat equation on a specified domain
- Convolution kernel
 - Stochastic kernel, the transition function of a stochastic process
 - Transition kernel, a generalization of a stochastic kernel
- Pricing kernel, the stochastic discount factor used in mathematical finance
- [Positive-definite kernel](#), a generalization of a positive-definite matrix
- Kernel trick, in statistics
- Reproducing kernel Hilbert space

Natural sciences [edit]

- Seed, inside the nut of most plants, especially:

Recap: basis functions in linear models for classification

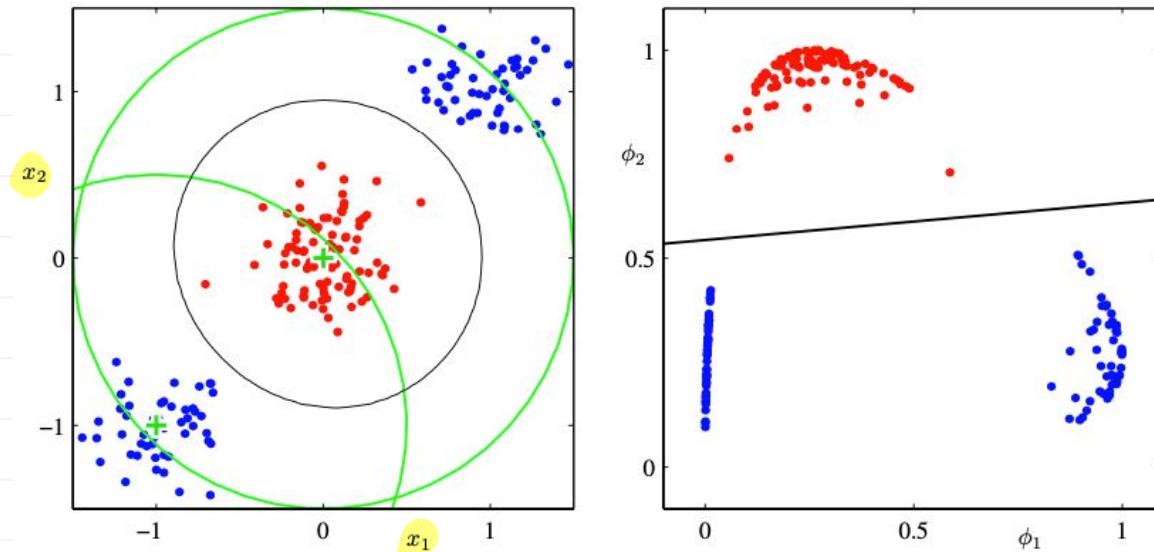
Linear w.r.t. Input x

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (4.4)$$

Linear w.r.t. basis $\phi(x)$

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

Fig 4.12



Recap: basis in linear models for regression

Linear w.r.t. input x

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

Linear w.r.t. basis $\phi(x)$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

fixed width

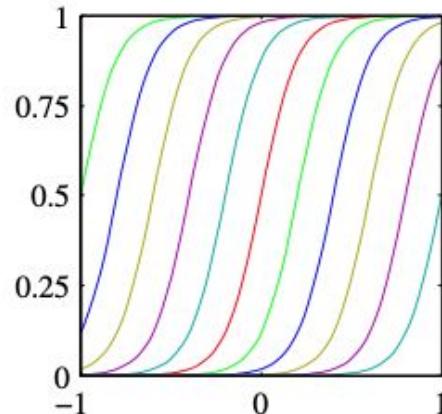
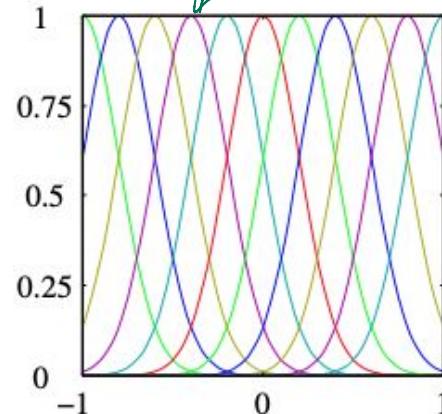
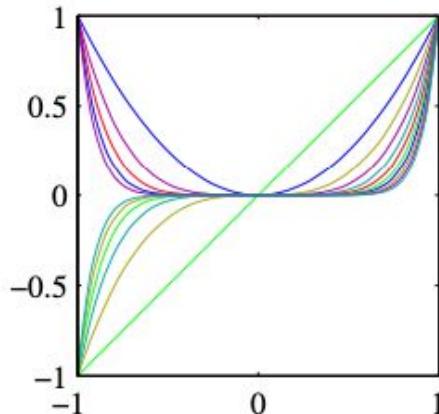


Figure 3.1 Examples of basis functions, showing polynomials on the left, Gaussians of the form (3.4) in the centre, and sigmoidal of the form (3.5) on the right.

So far

Basis functions CAN:

- Represent non-linear decision boundaries in the input space.
- Separate classes linearly (in the feature space) for classes that are not linearly separable in the input space.

Basis functions CAN NOT:

- Remove class overlap that already exist in input space. ← 
- Adapt its own shape to data (but neural nets can, “flexible basis functions”)
- Adapt the number of features to data (this lecture and next lecture ...)

$$x \in \mathbb{R}^d$$
$$w \in \mathbb{R}^m$$
$$y \sim w^T \phi(x)$$

"parametric"

Where are we going? Two key ideas

- Instead of “summarising” the training data into a set of weights (with fixed length), why not ask the features to adapt to data?
 - Continuity : Mostly targets don’t change abruptly.
 - Similarity : Each training pair (input, target) tells us something about the possible targets in the neighbourhood of the input.

Kernels formalise these ideas

- Nonparametric methods: do not rely on a fixed number of parameters, but rather usually on storing the entire training set (various loose definitions are used here).

“no parameters”
OR · unbounded # of “parameters”

How to get there?

- 0) Histograms
- 1) Simple density estimation kernels (must only be non-negative integrable)
- 2) Implicit feature mapping kernels (must be positive definite)

Warning:

- The term kernel is highly overloaded.
- Even today we consider two different types of kernel:
 - ➊ Smoothing kernel / density estimation kernel / Parzen window estimator kernel / Nadaraya-Watson kernel
 - ➋ Positive semidefinite kernel / reproducing kernel hilbert space kernel / implicit feature map kernel / support vector machine kernel / \approx Gaussian process covariance function or kernel

Density estimation

Observe data points $\{x_n\}_{n=1}^N$, draw i.i.d. from $p(x)$

Goal: estimate $p(x)$ from data

Figure 2.24 An illustration of the histogram approach to density estimation, in which a data set of 50 data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (2.241), with a common bin width Δ are shown for various values of Δ .

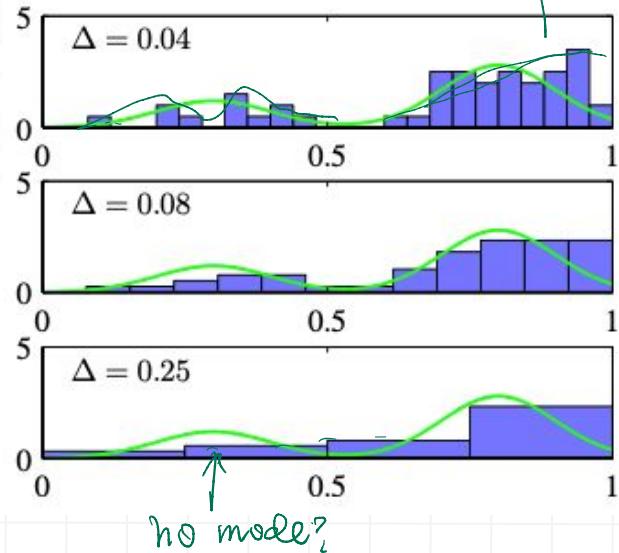
Histograms:

- Partition space into bins of width Δ
- Count the number of points in each bin
- Normalise

$$p_i = \frac{n_i}{N\Delta_i}$$

(2.241)

"Supervised ML" $\{(X, y)\}_{n=1}^N$
 $\{X_i\}_{i=1}^N$
 $\rightarrow p(x)$



Histograms as density estimators

Pros

- Each data point is used once. Applicable to sequentially arriving data.
- Data points can be thrown away once the histogram is computed -- fixed storage cost given bin width.
- Good for quickly visualising densities in 1 or 2 dimensions

what does a 2-dim histogram look like?

Cons

- Depend on bin width Δ
- Has discontinuities due to bin edges
- Does not work well in >2 dimensions - curse of dimensionality

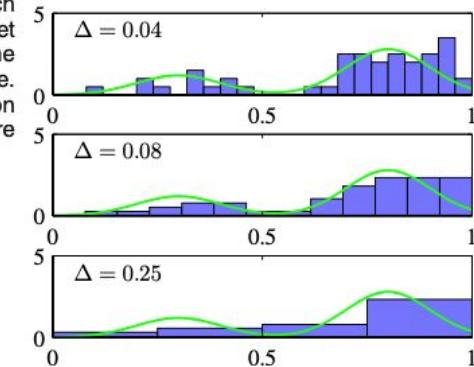
Implicit assumption: a distance measure in some local neighbourhood

Figure 2.24 An illustration of the histogram approach to density estimation, in which a data set of 50 data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (2.241), with a common bin width Δ are shown for various values of Δ .



Kernel density estimators

Nearest-neighbour methods



→ very few ~0 parameters

Non-parametric density estimation

N data points drawn i.i.d. from unknown distribution $p(x)$ in \mathbb{R}^D

Probability mass in a small region \mathcal{R}

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x}. \quad (2.242)$$

Prob that K points fall within \mathcal{R}

$$\text{Bin}(K|N, P) = \frac{N!}{K!(N-K)!} P^K (1-P)^{N-K}. \quad (2.243)$$

For large N

- Volume V of \mathcal{R} sufficiently small, s.t.
 $p(\mathbf{x})$ approximately constant in \mathcal{R}
- Volume V of \mathcal{R} sufficiently large, s.t.

K is sufficiently large and binomial sharply peaked

expected. $K \simeq NP$

$$P \simeq p(\mathbf{x})V \quad \text{volume of } \mathcal{R} \quad (2.245)$$

$$p(\mathbf{x}) = \frac{K}{NV}. \quad (2.246)$$

→ Fix K - K-nearest neighbours

Fix V - kernel density estimation

density est
classification
K means clustering

It can be shown that both the K-nearest-neighbour density estimator and the kernel density estimator converge to the true probability density in the limit $N \rightarrow \infty$ provided V shrinks suitably with N, and K grows with N (Duda and Hart, 1973).

Kernel density estimation (KDE), aka Parzen windows

$$p(\mathbf{x}) = \frac{K}{NV}. \quad (2.246)$$

Hypercubes

$$k(\mathbf{u}) = \begin{cases} 1, & |u_i| \leq 1/2, \\ 0, & \text{otherwise} \end{cases} \quad i = 1, \dots, D, \quad (2.247)$$

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right). \quad \text{how many points are } \pm \frac{1}{2} \text{ within } \mathbf{x}_n \quad (2.248)$$

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right) \quad (2.249)$$

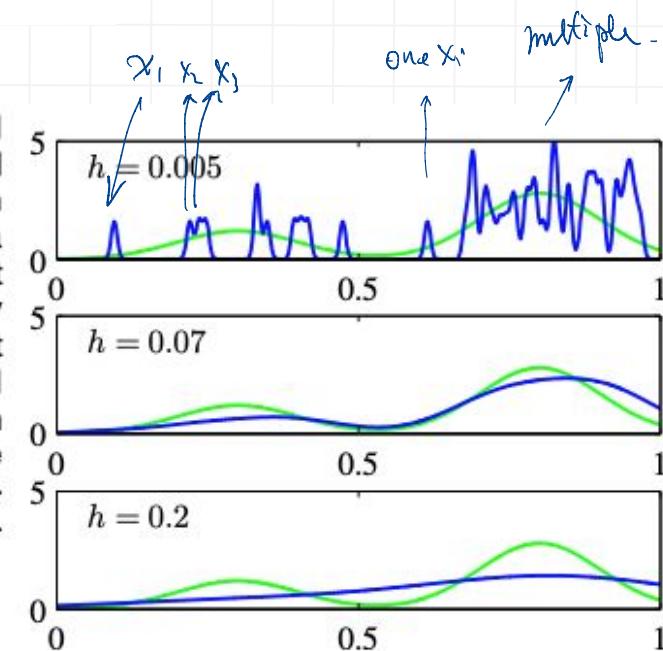
! We haven't improved from histograms (yet)

Kernel density estimation with Gaussians

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{1/2}} \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2} \right\} \quad (2.250)$$

Figure 2.25

Illustration of the kernel density model (2.250) applied to the same data set used to demonstrate the histogram approach in Figure 2.24. We see that h acts as a smoothing parameter and that if it is set too small (top panel), the result is a very noisy density model, whereas if it is set too large (bottom panel), then the bimodal nature of the underlying distribution from which the data is generated (shown by the green curve) is washed out. The best density model is obtained for some intermediate value of h (middle panel).



Kernel density estimation in general

Choose any $k(\mathbf{u})$ s.t.

$$k(\mathbf{u}) \geq 0, \quad (2.251)$$

$$\int k(\mathbf{u}) d\mathbf{u} = 1 \quad (2.252)$$

Pro: training data simply stored, no computation needed for “training”

Con: training data simply stored, expensive to evaluate the density

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right) \quad (2.249)$$

Nearest neighbour methods

$$p(\mathbf{x}) = \frac{K}{NV}.$$

(2.246)

Drawback of KDE: fixed V (not adapting to data), fixed $k(u)$, having to choose h (1-d kernel width parameter)

Alternative: fix K , use data to find V

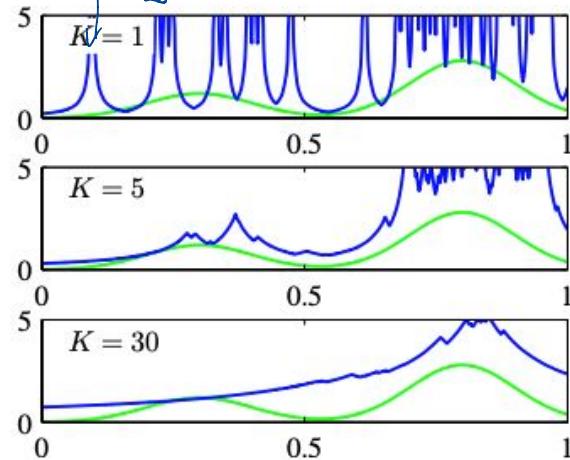
Consider a small sphere around x and then allow the radius to increase until it contains exactly K data points. Set V to the volume of the resulting sphere.

! Do not produce true density estimates.

Figure 2.26

Illustration of K -nearest-neighbour density estimation using the same data set as in Figures 2.25 and 2.24. We see that the parameter K governs the degree of smoothing, so that a small value of K leads to a very noisy density model (top panel), whereas a large value (bottom panel) smoothes out the bimodal nature of the true distribution (shown by the green curve) from which the data set was generated.

$V \rightarrow 0$ when $X \rightarrow X_n, n=1, \dots, N$ $p(x) \rightarrow \infty$



Nearest neighbour for classification

$$p(\mathbf{x}) = \frac{K}{NV}. \quad (2.246)$$

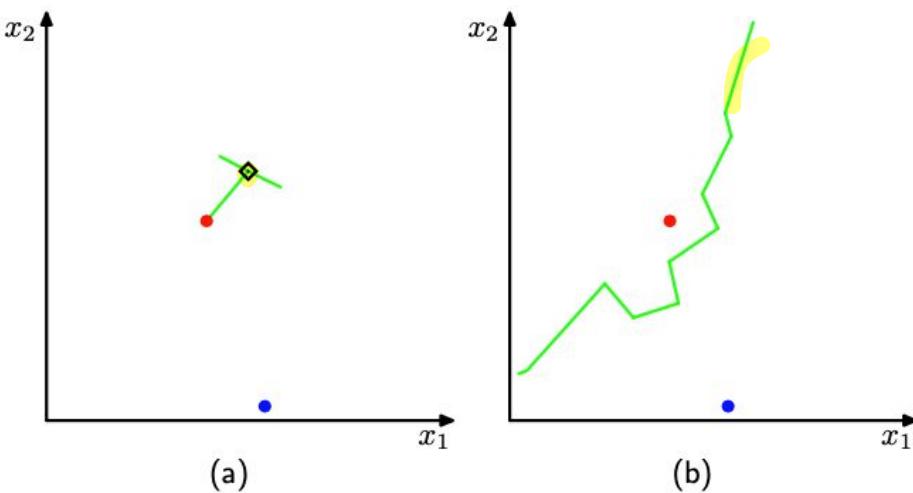
draw a sphere centred on \mathbf{x} containing precisely K points irrespective of their class.

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{K_k}{N_k V}. \quad (2.253)$$

$K = 1$: the *nearest-neighbour* rule

in the limit $N \rightarrow \infty$, the error rate of 1-NN is never more than twice the minimum achievable error rate of an optimal classifier, i.e., one that uses the true class distributions (Cover and Hart, 1967).

Figure 2.27 (a) In the K -nearest-neighbour classifier, a new point, shown by the black diamond, is classified according to the majority class membership of the K closest training data points, in this case $K = 3$. (b) In the nearest-neighbour ($K = 1$) approach to classification, the resulting decision boundary is composed of hyperplanes that form perpendicular bisectors of pairs of points from different classes.



larger $K \Rightarrow$ better classification

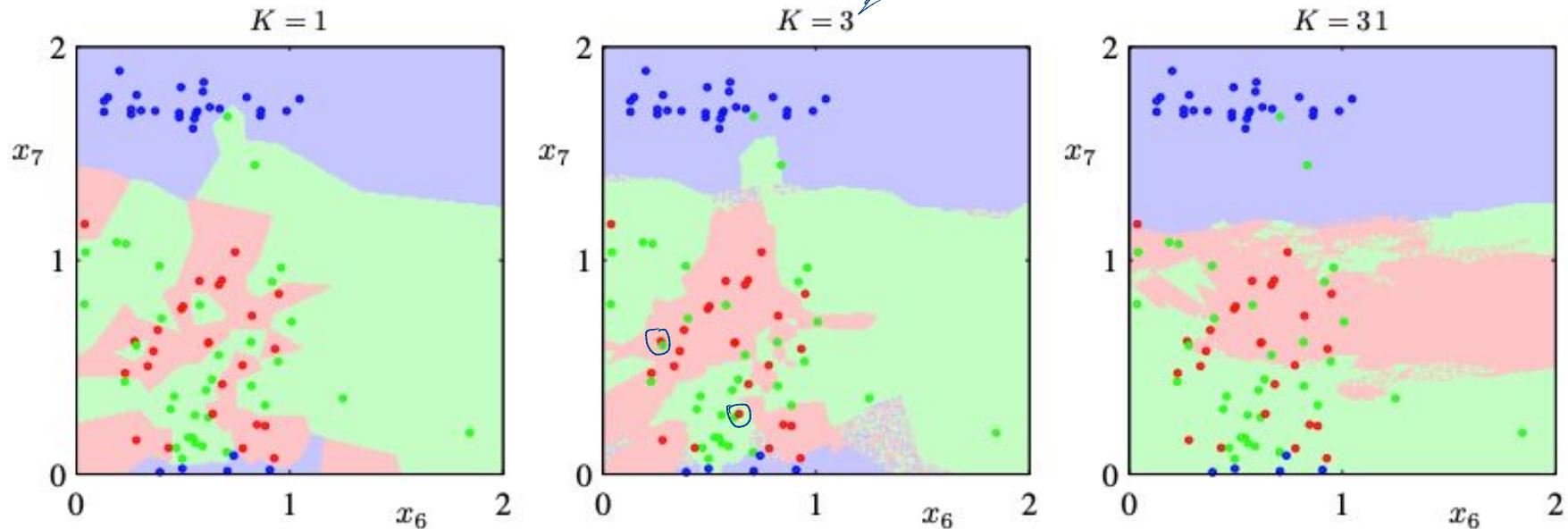


Figure 2.28 Plot of 200 data points from the oil data set showing values of x_6 plotted against x_7 , where the red, green, and blue points correspond to the 'laminar', 'annular', and 'homogeneous' classes, respectively. Also shown are the classifications of the input space given by the K -nearest-neighbour algorithm for various values of K .

Kernels

Basis functions recap

Non-parametric methods (Chap 2.5)

Dual representation (of linear models, Chap 6.1)

Constructing kernels (Chap 6.2)

The role of training data

- **Parametric methods**

$$w^T x, w^T \phi(x)$$

- Learn the model parameter w from the training data t .
- Discard the training data t .

- **Nonparametric methods**

- Use training data directly for prediction
- k -nearest neighbours : use k -closest data from the 'training' set for classification

circa 2005 ~ 2015
"non-parametric Bayesian"
 $\# \text{params} \rightarrow \infty$,

- **Kernel methods**

- Base prediction on linear combination of kernel functions evaluated at the training data.

Input vs Features

- A feature is a measurable property of a phenomenon being observed or any derived property thereof
 - raw features: the original data → Input
 - derived features: mappings of the original features to some other space (possibly high- or infinite dimensional, e.g., basis functions)
- Feature selection: which features matter for the problem at hand?
 - redundant features
 - problem dependent

w.r.t each other
are they informative of the target ?
- Feature extraction: can we combine the important features to a smaller set of new features?
 - compact representation versus ability to explain to a human

$$I(X_i; t) = 0$$

Kernel methods in one slide

- Consider a labelled training set $\{\mathbf{x}_i, t_i\}_{i=1}^N$ drawn i.i.d.
- On a new point \mathbf{x} , we will predict over all training points.

$$y(\mathbf{x}) = \sum_{i=1}^N a_i \cdot K(\mathbf{x}, \mathbf{x}_i)$$

"centered" on one training point \mathbf{x}_i

where $\{a_i\}_{i=1}^N$ are weights to be determined based on our training set, and $K(\cdot, \cdot)$ is a kernel function

- This is a major departure from the linear models considered previously!
- The kernel function measures the similarity between any two examples

- Prediction is a weighted average of the training targets
- Weights depend on the similarity of \mathbf{x} to each training example

$$y = \omega^\top \phi(\mathbf{x})$$

loosely
 $y(\mathbf{x}) \approx \sum_i \alpha_i \text{sim}(\mathbf{x}, \mathbf{x}_i)$

??

From feature function to kernels

MSE + L2 regularisation -

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \underline{\phi(\mathbf{x}_n)} - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (6.2) \quad \text{also (3.27)}$$

$$\Phi = \begin{matrix} \Phi \\ \scriptstyle N \times M \end{matrix} = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{bmatrix}$$

\leftarrow not what we need

$$\text{Optimal regularised w} \quad \mathbf{w}^* = (\lambda \mathbf{I} + \underline{\Phi^T \Phi})^{-1} \underline{\Phi^T \mathbf{t}} \quad \leftarrow \text{from chap 3}$$

Prediction function for new \mathbf{x}

$$y(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}^* = \phi(\mathbf{x})^T (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Goal \rightarrow

$y(\mathbf{x}) = \sum_{i=1}^N a_i \cdot K(\mathbf{x}, \mathbf{x}_i)$

$N \times N \text{ Gram matrix } \mathbf{K} = \underline{\Phi \Phi^T}$

$K \sim \text{sim}(x_i, x_j)$

$$\phi \phi^T \in \mathbb{R}^{N \times n}$$

$$\phi^T \phi \in \mathbb{R}^{M \times M}$$

Dual representation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} \quad (3.27)$$

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

$$\frac{\partial J}{\partial \mathbf{w}} = 0$$

$$(\Phi^\top \Phi + \lambda \mathbf{I}) \mathbf{w} = \Phi^\top \mathbf{t}$$

$$\lambda \mathbf{w} = \Phi^\top (\mathbf{t} - \Phi \mathbf{w})$$

$$\mathbf{w} = \Phi^\top \mathbf{a}$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n) a_n$$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^\top \underbrace{\Phi \Phi^\top}_{\text{Gram matrix } \mathbf{K}} \underbrace{\Phi \Phi^\top}_{\text{Gram matrix } \mathbf{K}} \mathbf{a} - \mathbf{a}^\top \underbrace{\Phi \Phi^\top}_{\text{Gram matrix } \mathbf{K}} \mathbf{t} + \frac{1}{2} \mathbf{t}^\top \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^\top \underbrace{\Phi \Phi^\top}_{\text{Gram matrix } \mathbf{K}} \mathbf{a}$$

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{bmatrix}$$

$N \times N$ Gram matrix $\mathbf{K} = \Phi \Phi^\top$

not $\Phi^\top \Phi$

$K^\top = K$

$$\mathbf{w}^* = (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$$

$$\overrightarrow{\mathbf{a}} \in \mathbb{R}^N \quad \overrightarrow{\mathbf{w}} \in \mathbb{R}^m$$

$$a_n = -\frac{1}{\lambda} \underbrace{\{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}}_{\text{pred. error on } x_n}$$

pred. error on x_n

Dual representation (contd)

$$J(\mathbf{a}) = \frac{1}{2}\mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2}\mathbf{t}^T \mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T \Phi \Phi^T \mathbf{a} \quad (6.5)$$

Define the $N \times N$ Gram matrix $\mathbf{K} = \underline{\Phi \Phi^T}$ with elements

$$\underline{K_{nm}} = \underline{\phi(\mathbf{x}_n)^T} \underline{\phi(\mathbf{x}_m)} = \underline{k(\mathbf{x}_n, \mathbf{x}_m)} \quad (6.6)$$

for x not in $\{\mathbf{x}_i\}_{i=1}^N$

$$\underline{k(x, x')} = \underline{\phi(x)^T} \underline{\phi(x')} = \sum_{i=1}^M \phi_i(x) \phi_i(x') \quad (6.10)$$

feat
of training
point

$$\underline{J(\mathbf{a})} = \frac{1}{2}\mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2}\mathbf{t}^T \mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T \mathbf{K} \mathbf{a}. \quad (6.7)$$

linear in \mathbf{a} linear in \mathbf{K} quadratic in \mathbf{a}

Solving for \mathbf{a} , prediction function

$$J(\mathbf{a}) = \frac{1}{2}\mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2}\mathbf{t}^T \mathbf{t} + \frac{\lambda}{2}\mathbf{a}^T \mathbf{K} \mathbf{a}. \quad (6.7)$$

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.$$

(6.8)

$$\mathbf{w} = \Phi^T \mathbf{a}$$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad (6.9)$$

Annotations:

- \mathbf{x} w. add $\{\mathbf{x}_n\}_{n=1}^N$ (pair-wise sim of training $\{\mathbf{x}_n\}_{n=1}^N$)
- sim. of \mathbf{x} w. add $\{\mathbf{x}_n\}_{n=1}^N$
- pre-compute during training \mathbf{R}^N (Inverse of $\mathbf{R}^{N \times N}$)
- $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$.

Prediction function for new \mathbf{x}

$$y(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}^* = \phi(\mathbf{x})^T (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

training points summarised in $(\mathbf{X}^T \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$

- So what is $K(\mathbf{x}, \mathbf{x}_i)$
- ① explicitly work it out using $\phi(\mathbf{x})$
 - ② $\phi(\mathbf{x})$ is implicit
 $K(\mathbf{x}_i, \mathbf{x}_j) \downarrow$ implies $\phi(\mathbf{x})$

GOAL

$$y(\mathbf{x}) = \sum_{i=1}^N a_i \cdot K(\mathbf{x}, \mathbf{x}_i)$$

Kernels

Basis functions recap

Non-parametric methods (Chap 2.5)

Dual representation (of linear models, Chap 6.1)

Constructing kernels (Chap 6.2)

Loosely, $K(x, x')$: similarity

RHS.

symmetric

$\forall \{x_n\}_{n=1}^N$

K positive semi definite.

The kernel function

- The **kernel function** is defined over two points, x and x' , of the input space

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x') \quad (6.10)$$

- $k(x, x')$ is symmetric.
- It is an inner product of two vectors of basis functions

$$k(x, x') = \langle \phi(x), \phi(x') \rangle.$$

- For prediction, the **kernel function** will be evaluated at the training data points. (See next slides.)

An example

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x') \quad (6.10)$$

further proof.

$$k(\mathbf{x}, \mathbf{z}) = (\underline{\mathbf{x}^T \mathbf{z}})^2. \quad (\underline{x_1, x_2}) \quad (\underline{z_1, z_2}) \quad (6.11)$$

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\ &= \underline{\phi(\mathbf{x})^T \phi(\mathbf{z})}. \end{aligned}$$

*only involve
(x₁, x₂)*

*only involve
(z₁, z₂)*

$$\phi''(\mathbf{x}) = (x_1, x_2, x_1 x_2, x_1^2, x_2^2)^T$$

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$$

$$\underline{\phi(\mathbf{x})} = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$$

*features of polynomial
order 2*

More generally, however, we need a simple way to test whether a function constitutes a valid kernel without having to construct the function $\phi(\mathbf{x})$ explicitly. A **necessary and sufficient condition for a function $k(\mathbf{x}, \mathbf{x}')$ to be a valid kernel** (Shawe-Taylor and Cristianini, 2004) is that the **Gram matrix \mathbf{K}** , whose elements are given by $k(\mathbf{x}_n, \mathbf{x}_m)$, should be **positive semidefinite for all possible choices of the set $\{\mathbf{x}_n\}$** . Note that a **positive semidefinite matrix** is not the **same thing as a matrix whose elements are nonnegative**.

all eigenvalues ≥ 0

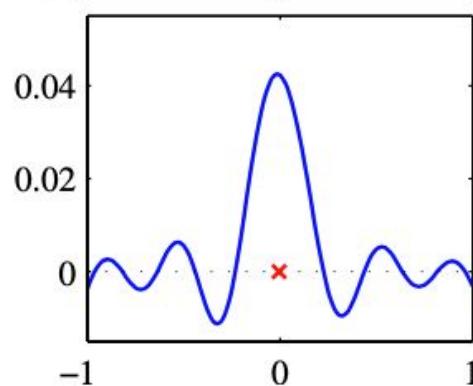
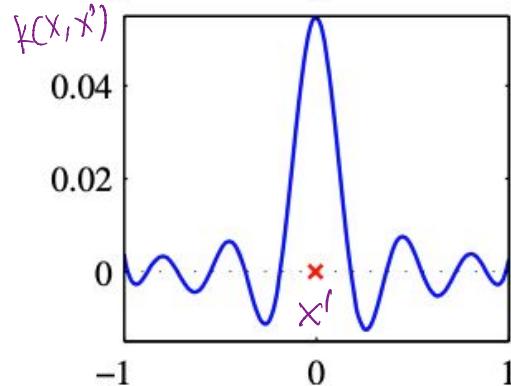
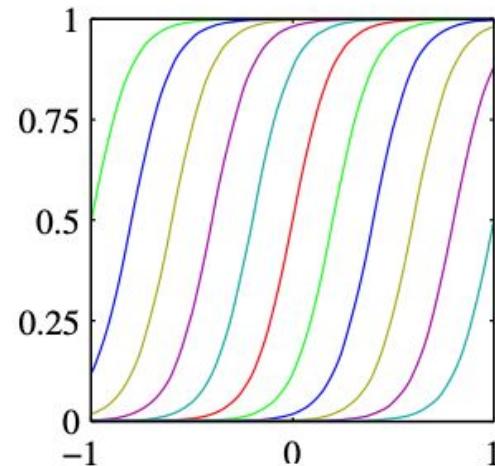
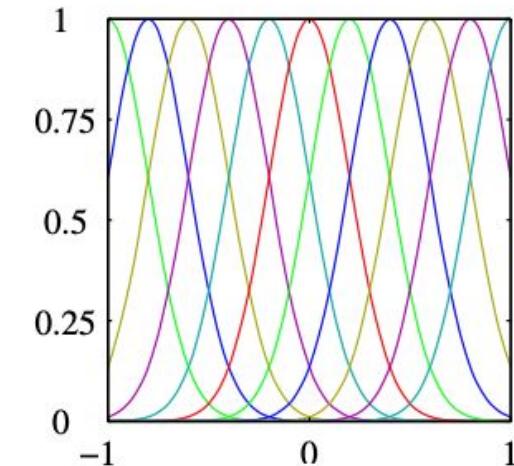
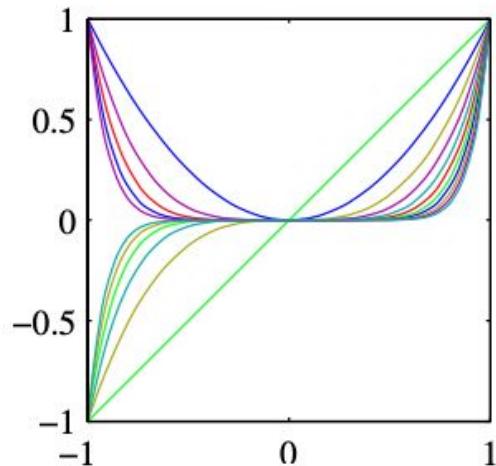


Figure 6.1 Illustration of the construction of kernel functions starting from a corresponding set of basis functions. In each column the lower plot shows the kernel function $k(x, x')$ defined by (6.10) plotted as a function of x for $x' = 0$, while the upper plot shows the corresponding basis functions given by polynomials (left column), ‘Gaussians’ (centre column), and logistic sigmoids (right column).

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\mathbf{x})$ is a function from \mathbf{x} to \mathbb{R}^M , $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M , \mathbf{A} is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

$$\mathbf{x}^T \mathbf{x}'$$

$$(\mathbf{x}^T \mathbf{x}')^2$$

Gaussian / Radial Basis Function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma^2) \quad (6.23) \quad \rightarrow \text{What's } \phi(\mathbf{x})?$$

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + (\mathbf{x}')^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}' \quad (6.24)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\mathbf{x}^T \mathbf{x}/2\sigma^2) \exp(\mathbf{x}^T \mathbf{x}'/\sigma^2) \exp(-(\mathbf{x}')^T \mathbf{x}'/2\sigma^2) \quad (6.25)$$

$$\begin{array}{ccc} f(x) & \exp\left(\frac{-1}{\sigma^2} k(x, x')\right) & f(x') \\ & \downarrow & \\ & x^T x' & \end{array}$$

$$\phi(x) = x$$

Further examples of kernels

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^M$$

only terms of degree M

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^M$$

all terms up to degree M

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$$

Gaussian kernel

$$k(\mathbf{x}, \mathbf{x}') = \tanh(a \mathbf{x}^\top \mathbf{x}' + b)$$

Sigmoidal kernel (invalid)

easier to compute
that $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$

Generally, we call

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$$

only depend
vec. diff

linear kernel

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$$

stationary kernel

$$k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$$

homogeneous kernel

only depend on distance

Kernels over graphs, strings, sets

- We 'only' need an appropriate similarity measure $k(\mathbf{x}, \mathbf{x}')$ which is a kernel.
- Example: Given a set \mathcal{A} and the set of all subsets of \mathcal{A} , called the **power set** $\mathcal{P}(\mathcal{A})$.
- For two subsets $\mathcal{A}_1, \mathcal{A}_2 \in \mathcal{P}(\mathcal{A})$, denote the number of elements of the intersection of \mathcal{A}_1 and \mathcal{A}_2 by $|\mathcal{A}_1 \cap \mathcal{A}_2|$.
- Then it can be shown that

$$k(\mathcal{A}_1, \mathcal{A}_2) = 2^{|\mathcal{A}_1 \cap \mathcal{A}_2|}$$

corresponds to an inner product in a feature space.
Therefore, $k(\mathcal{A}_1, \mathcal{A}_2)$ is a valid kernel function.

Kernels from probabilistic generative models

- Given $p(\mathbf{x})$, we can define a kernel

$$k(\mathbf{x}, \mathbf{x}') = p(\mathbf{x}) p(\mathbf{x}'),$$

which means two inputs \mathbf{x} and \mathbf{x}' are similar if they both have high probabilities.

- Include a weighting function $p(i)$ and extend the kernel to

$$k(\mathbf{x}, \mathbf{x}') = \sum_i p(\mathbf{x} | i) p(\mathbf{x}' | i) p(i).$$

- For a continuous variable \mathbf{z}

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{x}' | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}.$$

- Hidden Markov Model with sequences of length L .

Kernels for regression and classification: summary

- Pick a suitable kernel function $k(\mathbf{x}, \mathbf{x}')$
 - e.g. by computing inner product of some basis functions
- Make predictions by suitably combining $k(\mathbf{x}, \mathbf{x}_n)$ for each training example \mathbf{x}_n
 - implicitly, a linear model in some high-dimensional space
- For linear regression, we go from

$$y(\mathbf{x}) = \phi(\mathbf{x})^\top (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$$

to

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

- can plug in suitable kernel function to implicitly perform nonlinear transformation

no need to write out $\phi(\mathbf{x})$

Cons of $k(\mathbf{x}, \mathbf{x}')$
 $\mathcal{O}(N^2)$
 $(\mathbf{K} + \lambda \mathbf{I}_N)^{-1}$ $\mathcal{O}(N^3)$

benefits of using $k(\mathbf{x}, \mathbf{x}')$

- $\phi(\mathbf{x})$ can be implicitly "infinite" dimensions
- Compose new \mathbf{k} from existing \mathbf{k}
- "sparse" \mathbf{x}
- RBF SVMs \rightarrow "sparse" weighted nearest neighbour
- New toolset

Kernels

Basis functions recap

Non-parametric methods (Chap 2.5)

Dual representation (of linear models, Chap 6.1)

Constructing kernels (Chap 6.2)