

Lecture 10a

Kernels

Outline

Making ML models nonlinear

- ▶ Feature Expansions

From Feature Expansions to Kernels

- ▶ Computational benefits
- ▶ Kernelizing the distance computation
- ▶ Kernelizing the one-class SVM

From Kernels to Feature Expansions

- ▶ Testing positive semi-definiteness
- ▶ Commonly used kernels

From Linear to Nonlinear Models

Observation:

- ▶ Linear or quadratic models (e.g. Fisher discriminant, linear regression, etc.) are easy to train, but often lack the representation power to learn complex tasks.

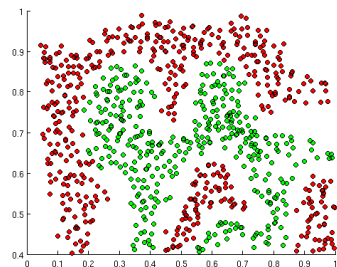


Image source: Ng, Stanford OpenClassroom ML course

Idea:

- ▶ No need to invent new algorithms. Apply some nonlinear transformation $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^h$ to the data, and learn a *linear* model in the image domain:

$$f(\mathbf{x}) = \boldsymbol{\beta}^\top \phi(\mathbf{x}) + b$$

- ▶ The function ϕ is usually called a ‘feature map’.

Part 1

Feature Expansions

Quadratic Features Expansions

General expression:

$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ (x_i)_{i=1}^d \\ (x_i x_j)_{i,j=1}^d \end{bmatrix} \in \mathbb{R}^{(d+1) \cdot (d+2)/2}$$

Observation:

- ▶ If building a linear model on top of ϕ , we can show that it is equivalent to a quadratic model of \mathbf{x} :

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$$

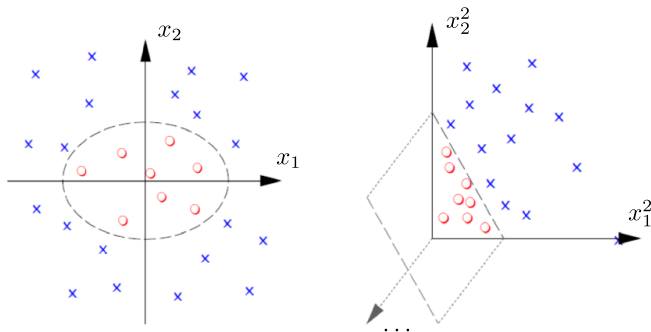
for some parameters A, \mathbf{b}, c .

- ▶ We can use existing methods (e.g. PCA, linear regression, Fisher discriminant, ...), and increase their predictive power.

Linear Model in Feature Space

Example for $d = 2$:

$$\begin{aligned}\phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^6 \\ (x_1, x_2) &\mapsto (1, x_1, x_2, x_1^2, x_2^2, x_1x_2)\end{aligned}$$



Observation:

- The nonlinear (quadratic) problem is made linear in feature space.

From simple to more powerful discriminants

Difference of means:

- finds a linear projection that maximizes the difference between the two classes.

$$f(\mathbf{x}) = \mathbf{x}^\top (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$$

Fisher discriminant:

- finds a linear projection that maximizes the separability of the two classes.

$$f(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$$

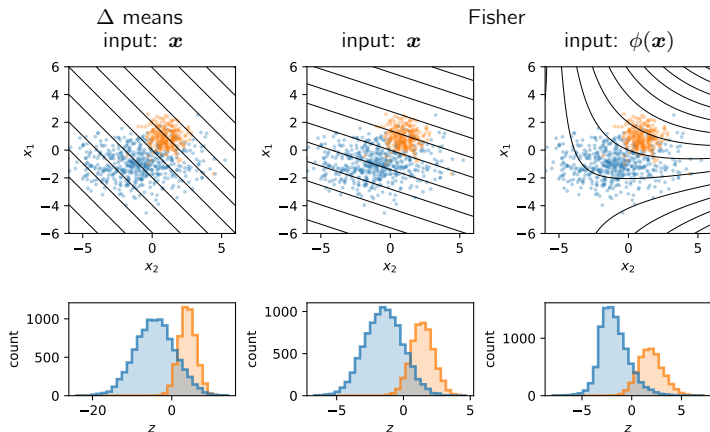
Fisher discriminant (with quadratic features):

- finds a nonlinear mapping that maximizes the separability of the two classes.

$$f(\mathbf{x}) = \phi(\mathbf{x})^\top \Sigma_\phi^{-1} (\boldsymbol{\mu}_2^\phi - \boldsymbol{\mu}_1^\phi)$$

where $\boldsymbol{\mu}_1^\phi$, $\boldsymbol{\mu}_2^\phi$, Σ_ϕ indicate the means and covariances computed on $\phi(\mathbf{x})$ instead of \mathbf{x} .

Example



Observations:

- The nonlinear transformation gives more flexibility for increasing class separability, e.g. by bending the contour lines to mitigate the higher variance in the first class.

Other Nonlinear Expansions

Polynomial expansions:

- Generalization of quadratic expansions to higher degree:

$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ (x_i)_{i=1}^d \\ (x_i x_j)_{i,j=1}^d \\ (x_i x_j x_k)_{i,j,k=1}^d \\ \vdots \end{bmatrix}$$

Hashing:

- Map to binary vectors testing for membership to some subspace, e.g.

$$\phi(\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \sigma(\mathbf{w}_2^\top \mathbf{x} + b_2) \\ \vdots \\ \sigma(\mathbf{w}_h^\top \mathbf{x} + b_h) \end{bmatrix}$$

with $(\mathbf{w}_i, b_i) \sim p(\mathbf{w}, b)$ some distribution and σ is some nonlinear function, e.g. a sign function, or a sinusoid.

Part 2

From Feature Expansions to Kernels

Distance Computation in Feature Space

Many prediction functions can be expressed in terms of dot products between data points in feature space (called 'kernels'):

Example:

- Distance to the mean:

$$\begin{aligned}d(\mathbf{x}) &= \left\| \phi(\mathbf{x}) - \overbrace{\frac{1}{N} \sum_i \phi(\mathbf{x}_i)}^{\boldsymbol{\mu}^\phi} \right\| \\&= \left(\left\| \phi(\mathbf{x}) - \frac{1}{N} \sum_i \phi(\mathbf{x}_i) \right\|^2 \right)^{\frac{1}{2}} \\&= \left(\underbrace{\phi(\mathbf{x})^\top \phi(\mathbf{x})}_{\text{dot product}} - \frac{2}{N} \sum_i \underbrace{\phi(\mathbf{x})^\top \phi(\mathbf{x}_i)}_{\text{dot product}} + \frac{1}{N^2} \sum_{i,j} \underbrace{\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)}_{\text{dot product}} \right)^{\frac{1}{2}}\end{aligned}$$

Question: Can we bypass the computation of the feature map when evaluating dot products?

The Kernel Trick

Example: Dot product computations can be computed very quickly for certain choices of feature maps ϕ (we refer to such dot products as kernels)

$$\begin{aligned} \left\langle \overbrace{\begin{bmatrix} 1 \\ \sqrt{2}(x_i)_{i=1}^d \\ (x_i x_j)_{i,j=1}^d \end{bmatrix}}^{\phi(\mathbf{x})}, \overbrace{\begin{bmatrix} 1 \\ \sqrt{2}(x'_i)_{i=1}^d \\ (x'_i x'_j)_{i,j=1}^d \end{bmatrix}}^{\phi(\mathbf{x}')}\right\rangle &= 1 + \sqrt{2}\sqrt{2} \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j \\ &= 1 + 2 \sum_{i=1}^d x_i x'_i + \left(\sum_{i=1}^d x_i x'_i \right) \left(\sum_{j=1}^d x_j x'_j \right) \\ &= \left(1 + \sum_{i=1}^d x_i x'_i \right)^2 \\ &= \underbrace{\left(1 + \langle \mathbf{x}, \mathbf{x}' \rangle \right)^2}_{k(\mathbf{x}, \mathbf{x}')} \end{aligned}$$

Observation:

- ▶ Evaluating $k(\mathbf{x}, \mathbf{x}')$ only requires $\approx d$ computations (instead of $\approx d^2$ computations if evaluating ϕ).

Computing ϕ vs. Computing k

Cost of making one prediction:

$$f(\mathbf{x}) = \underbrace{\sum_i \alpha_i \phi(\mathbf{x}_i)}_{\mathbf{w}}^\top \underbrace{\phi(\mathbf{x})}_{k(\mathbf{x}_i, \mathbf{x})} + \theta$$

→ approach 1 (classic): $O(h)$

→ approach 2 (kernel): $O(N \cdot d)$

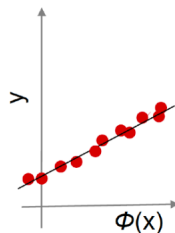
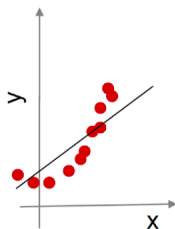
Examples:

N	d	(ϕ, k)	Approach 1 (classic)	Approach 2 (kernel)
1000	25	Linear	25	—
1000	25	Poly($q = 2$)	625	25000
1000	25	Poly($q = 3$)	15625	25000
1000	25	Poly($q = 4$)	390625	25000
1000	25	Poly($q = 5$)	9765625	25000

Kernel Ridge Regression

Ideas:

- ▶ Perform standard ridge regression in some feature space.



- ▶ Rewrite the resulting model in a way that the feature map computations only appear as dot products (so that they can be replaced by kernels).

Kernel Ridge Regression

- ▶ Redefine the prediction function $y = \mathbf{w}^\top \phi(\mathbf{x})$ where $\mathbf{w} \in \mathbb{R}^h$ and the objective to minimize as

$$\mathcal{E}(\mathbf{w}) = \frac{1}{N} \sum_{k=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_k) - t_k)^2 \quad \text{subject to} \quad \|\mathbf{w}\|^2 \leq C.$$

- ▶ As we have seen in the lecture on regression, the solution to this optimization problem has the form:

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

where $\phi(X) = (\phi(\mathbf{x}_1) | \dots | \phi(\mathbf{x}_N))$, and for an appropriate choice of parameter λ .

- ▶ A new data point is then predicted as

$$\begin{aligned} y = \mathbf{w}^\top \phi(\mathbf{x}) &= \phi(\mathbf{x})^\top \mathbf{w} \\ &= \phi(\mathbf{x})^\top (\phi(X)\phi(X)^\top + \lambda I)^{-1} \phi(X)\mathbf{t} \end{aligned}$$

Kernel Ridge Regression

from: $y = \phi(x)^\top (\underbrace{\phi(X)\phi(X)^\top}_{\bar{\Phi}} + \lambda I)^{-1} \phi(X)t$

$$y = \phi(x)^\top \underbrace{(\bar{\Phi} + \lambda I)^{-1}}_{\substack{\bar{\Phi}\bar{\Phi}^\top\bar{\Phi} + \lambda\bar{\Phi} \\ (\bar{\Phi}\bar{\Phi}^\top + \lambda I)\bar{\Phi}}} \bar{\Phi} \underbrace{(\bar{\Phi}^\top\bar{\Phi} + \lambda I)^{-1}t}_{\bar{\Phi}^\top\bar{\Phi} + \lambda I}$$

$$= \underbrace{\phi(x)^\top}_{1 \times h} \underbrace{\bar{\Phi}}_{N \times N} \underbrace{(\bar{\Phi}^\top\bar{\Phi} + \lambda I)^{-1}}_{N \times 1} +$$

$$= k(x, X)(k(X, X) + \lambda I)^{-1} \cdot t$$

we arrive at:

$$y = k(x, X)(K + \lambda I)^{-1}t$$

Kernel Fisher Discriminant

We can kernelize the Fisher discriminant model by first observing that the projection on the Fisher discriminant can be written as:

$$\begin{aligned}y &= \phi(\mathbf{x})^\top \mathbf{w} \\&= \phi(\mathbf{x})^\top (\Sigma_\phi + \lambda I)^{-1} \cdot (\boldsymbol{\mu}_2^\phi - \boldsymbol{\mu}_1^\phi) \\&= \phi(\mathbf{x})^\top \left(\frac{1}{N} \phi(X) \phi(X)^\top + \lambda I \right)^{-1} (\boldsymbol{\mu}_2^\phi - \boldsymbol{\mu}_1^\phi) \\&= \phi(\mathbf{x})^\top \left(\frac{1}{N} \phi(X) \phi(X)^\top + \lambda I \right)^{-1} \phi(X) \mathbf{t}\end{aligned}$$

where

$$t_i = \begin{cases} \frac{1}{|\mathcal{G}_2|} & \text{if } i \in \mathcal{G}_2 \\ -\frac{1}{|\mathcal{G}_1|} & \text{if } i \in \mathcal{G}_1 \end{cases}$$

and then proceed similarly as for kernelizing ridge regression.

Further Kernelizations

Many linear models can be kernelized.

Examples:

- ▶ Kernel PCA
- ▶ Kernel CCA
- ▶ Kernel SVM
- ▶ Kernel SVR
- ▶ Kernel One-Class SVM

Often, the term 'kernel' is omitted (e.g. the SVM modules in scikit learn uses by default a nonlinear kernel).

Part 3

From Kernels to Feature Spaces

Can we Choose any Function k in Practice?

Negative example:

- ▶ Suppose you would like to use the candidate kernel $k(\mathbf{x}, \mathbf{x}') = I(\mathbf{x} \neq \mathbf{x}')$ in an algorithm based on distance computation.
- ▶ Similarly to what we did before, the distance can be kernelized as:

$$\begin{aligned} d(\mathbf{x}, \mathbf{x}') &= \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\| \\ &= \sqrt{\|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|^2} \\ &= \sqrt{k(\mathbf{x}, \mathbf{x}) - 2k(\mathbf{x}, \mathbf{x}') + k(\mathbf{x}', \mathbf{x}')}. \end{aligned}$$

Injecting our candidate kernel into this expression, we get

$$\begin{aligned} &= \sqrt{I(\mathbf{x} \neq \mathbf{x}) - 2I(\mathbf{x} \neq \mathbf{x}') + I(\mathbf{x}' \neq \mathbf{x}')} \\ &= \sqrt{0 - 2 + 0} \end{aligned}$$

which throws an error, because the square root does not apply to negative numbers.

Conclusion:

- ▶ Not all functions are valid kernels. Some conditions need to be verified.

Testing the Validity of a Kernel

[Mercer] If k is continuous, symmetric, and for all square-integrable functions $g(\mathbf{x})$ satisfies the condition:

$$\iint g(\mathbf{x}) k(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0$$

then the kernel can be expanded as:

$$k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{N_F} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{x}')$$

with $\lambda_i > 0$, and $N_F \in \mathbb{N}$ or $N_F = \infty$ and we can construct the feature map

$$\phi(\mathbf{x}) := \begin{pmatrix} \sqrt{\lambda_1} \psi_1(\mathbf{x}) \\ \sqrt{\lambda_2} \psi_2(\mathbf{x}) \\ \vdots \end{pmatrix}$$

satisfying $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$

Testing the Validity of a Kernel

In practice, to test whether the kernel is valid, we need to:

1. verify that the kernel is symmetric, i.e.:

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

for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$.

2. show that it is positive semi-definite, i.e.:

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

for all sequences of points $(\mathbf{x}_1, \dots, \mathbf{x}_N)$ and sequences of real numbers (c_1, \dots, c_N) .

Example

Show that $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ with $k(\mathbf{x}, \mathbf{x}') = \alpha + x_1 x'_1$ and $\alpha > 0$, is a Mercer kernel.

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

$$\begin{aligned} & \sum_{i=1}^N \sum_{j=1}^N c_i c_j (\alpha + x_1^{(i)} \cdot x_1^{(j)}) \\ &= \sum_i \sum_j c_i c_j \alpha + \sum_i \sum_j c_i c_j x_1^{(i)} x_1^{(j)} \\ &= \underbrace{\alpha \left(\sum_i c_i \right) \left(\sum_j c_j \right)}_{\geq 0} + \underbrace{\left(\sum_i c_i x_1^{(i)} \right) \left(\sum_j c_j x_1^{(j)} \right)}_{\geq 0} \geq 0 \end{aligned}$$

Find a feature map $\phi(\mathbf{x})$ associated to this kernel

$$\phi(\mathbf{x}) = \begin{pmatrix} \sqrt{\alpha} \\ x_1 \end{pmatrix} \quad \left(\begin{pmatrix} \sqrt{\alpha} \\ x_1 \end{pmatrix} \right)^T \left(\begin{pmatrix} \sqrt{\alpha} \\ x'_1 \end{pmatrix} \right) = \alpha + x_1 x'_1$$

Example of Kernels

Examples of commonly used kernels satisfying the Mercer property

Linear	$k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$	
Polynomial	$k(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + \beta)^\gamma$	$\beta \in \mathbb{R}_{\geq 0}, \gamma \in \mathbb{N}$
Gaussian	$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \ \mathbf{x} - \mathbf{x}'\ ^2)$	$\gamma \in \mathbb{R}_{> 0}$
Student	$k(\mathbf{x}, \mathbf{x}') = \frac{1}{\alpha + \ \mathbf{x} - \mathbf{x}'\ ^2}$	$\alpha \in \mathbb{R}_{> 0}$

Summary

Summary

- ▶ One approach to make models nonlinear is to pass the data through a nonlinear feature map ϕ , and then use a linear model on top of it.
- ▶ Feature maps can either be engineered based on existing knowledge, or it can be of generic type (e.g. quadratic expansion, polynomial expansion, or random hashing).
- ▶ Many learning algorithms applied in feature space can be expressed in terms of dot products in feature space.
- ▶ For certain types of feature maps, dot products can be computed without explicitly mapping to the feature space, and result in much faster computation, and we call these dot product computations 'kernels'.
- ▶ In practice, one cannot choose any kernel function. The kernel must be symmetric and positive semi-definite so that the existence of a corresponding feature map is ensured, and so that the various ML algorithms are guaranteed to remain applicable.