Kernel Methods

CS771: Introduction to Machine Learning Pivush Rai

Linear Models for Nonlinear Problems?

Consider the following one-dimensional inputs from two classes



Can't separate using a linear hyperplane



Linear Models for Nonlinear Problems?

• Consider mapping each x to two-dimensions as $x \rightarrow z = [z_1, z_2] = [x, x^2]$



Classes are now linearly separable in the two-dimensional space



Linear Models for Nonlinear Problems

Can assume a feature mapping ϕ that maps/transforms the inputs to a "nice" space



.. and then happily apply a linear model in the new space!



Not Every Mapping is Helpful

- Not every higher-dim mapping helps in learning nonlinear patterns
- Must be a <u>nonlinear</u> mapping
- For the nonlin classfn problem we saw earlier, consider some possible mappings





How to get these "good" (nonlinear) mappings?

Learn good mappings from data itself (e.g., deep learning or distance metric learning)

6

Use pre-defined "good" mappings (e.g., defined by kernel functions - today's topic)



Some Pre-defined Kernel Functions

• Linear kernel: $k(x, z) = x^{T}z$

Several other kernels proposed for non-vector data, such as trees, strings, etc

Remember that kernels are a notion of similarity between pairs of inputs

or can be learned from data (a bit

advanced for this course)



- Quadratic Kernel: $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^{\mathsf{T}} \mathbf{z})^2$ or $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^{\mathsf{T}} \mathbf{z})^2$
- Polynomial Kernel (of degree d): $k(x, z) = (x^{T}z)^{d}$ or $k(x, z) = (1 + x^{T}z)^{d}$
- Radial Basis Function (RBF) or "Gaussian" Kernel: $k(x, z) = \exp[-\gamma ||x z||^2]$
 - Gaussian kernel gives a similarity score between 0 and 1
 - $\gamma > 0$ is a hyperparameter (called the kernel bandwidth parameter)
 - The RBF kernel corresponds to an infinite dim. feature space \mathcal{F} (i.e., you can't actually write down or store the map $\phi(x)$ explicitly – but we don't need to do that anyway S)
 - Also called "stationary kernel": only depends on the distance between x and z (translating both by the same amount won't change the value of k(x, z)
- Which kernel to use or its hyperparams (e.g., d, γ) values can be set via cross-val. CS771: Intro to ML

Controls how the distance between two inputs should be converted into a similarity

Kernels as (Implicit) Feature Maps

- Consider two inputs (in the same two-dim feature space): $\mathbf{x} = [x_1, x_2], \mathbf{z} = [z_1, z_2]$
- Suppose we have a function k(.,.) which takes two inputs $m{x}$ and $m{z}$ and computes

Called the Can think of this as a notion "kernel function" $k(\mathbf{x},\mathbf{z}) = (\mathbf{x}^{\mathsf{T}}\mathbf{z})^2$ of similarity b/w \boldsymbol{x} and \boldsymbol{z} $= (x_1 z_1 + x_2 z_2)^2$ Didn't need to compute $\phi(x)$ explicitly. Just using the definition of the kernel k(x,z) = $= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2$ $(x^{\mathsf{T}}z)^2$ implicitly gave us this mapping for each input $= (x_1^2, \sqrt{2}x_1x_2, x_2^2)^{\top} (z_1^2, \sqrt{2}z_1z_2, z_2^2)$ Thus kernel function $k(\mathbf{x}, \mathbf{z}) =$ $(\mathbf{x}^{\mathsf{T}}\mathbf{z})^2$ implicitly defined a feature mapping $= \phi(\mathbf{x})^{\mathsf{T}} \phi(\mathbf{z})$ Dot product similarity in ϕ such that for $x = [x_1, x_2]$, the new feature space $\phi(\mathbf{x}) = \left(x_1^2, \sqrt{2}x_1x_2, x_2^2\right)$ defined by the mapping ϕ

This is not a dot/inner product similarity but similarity using a more general function of **x** and **z** (square of dot product)

> Remember that a kernel does two things: Maps the data implicitly into a new feature space (feature transformation) and computes pairwise similarity between any two inputs under the new feature representation



• Also didn't have to compute $\phi(x)^{\top}\phi(z)$. Defn $k(x,z) = (x^{\top}z)^2$ gives that

RBF Kernel = Infinite Dimensional Mapping

- We saw that the RBF/Gaussian kernel is defined as $k(x, z) = \exp[-\gamma ||x z||^2]$
- Using this kernel corresponds to mapping data to infinite dimensional space

$$k(x,z) = \exp[-(x-z)^{2}] \quad (\text{assuming } \gamma = 1 \text{ and } x \text{ and } z \text{ to be scalars})$$

$$= \exp(-x^{2}) \exp(-z^{2}) \exp(2xz)$$

$$= \exp(-x^{2}) \exp(-z^{2}) \sum_{0=1}^{\infty} \frac{2^{k} x^{k} z^{k}}{k!}$$

$$= \phi(x)^{T} \phi(z) \qquad \text{Thus an infinite-dim vector (ignoring the constants coming from the 2^{k} and k! terms}$$

- Here $\phi(x) = [\exp(-x^2)x^0, \exp(-x^2)x^1, \exp(-x^2)x^2, \exp(-x^2)x^3, \dots, \exp(-x^2)x^{\infty}]$
- But again, note that we never need to compute $\phi(x)$ to compute k(x,z)
 - k(x,z) is easily computable from its definition itself $(\exp[-(x-z)^2])$ in this case)

Kernel Function: Some Other Aspects

- Not every function of the form $k(x, z) = \phi(x)^{\mathsf{T}} \phi(z)$ is a kernel function
- k must satisfy Mercer's Condition
 - k must define a dot product for some Hilbert Space
 - Above is true if k is symmetric and positive semi-definite (p.s.d.) function (though there are exceptions; there are also "indefinite" kernels) Loosely speaking a PSD <u>function here</u>

For all "square integrable" functions
$$f$$

(such functions satisfy $\int f(x)^2 dx < \infty$
$$\iint f(x)k(x,z)f(z)dxdz \geq 0$$

Loosely speaking a PSD function here
means that if we evaluate this function for
$$N$$
 inputs (N^2 pairs) then the $N \times N$ matrix
will be PSD (also called a kernel matrix)

Can easily verify that the Mercer's Condition holds for these

CS771: Intro to ML

- Let k_1, k_2 be two kernel functions then the following are as well
 - $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$: simple sum
 - $k(\mathbf{x}, \mathbf{z}) = \alpha k_1(\mathbf{x}, \mathbf{z})$: scalar product with $\alpha > 0$
 - $k(x, z) = k_1(x, z)k_2(x, z)$: direct product of two kernels

Can also combine these rules and the resulting function will also be a kernel function

Kernel Matrix

- Kernel based ML algos work with kernel matrices rather than feature vectors
- Given N inputs, the kernel function k can be used to construct a Kernel Matrix K
- The kernel matrix **K** is of size $N \times N$ with each entry defined as

$$K_{ij} = k(\boldsymbol{x}_i, \boldsymbol{x}_j) = \phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}_j)$$

Note again that we don't need to compute ϕ and this dot product explicitly

• K_{ij} : Similarity between the i^{th} and j^{th} inputs in the kernel induced feature space ϕ



Using Kernels in ML algorithms



17

Using Kernels

- Kernels can turn many linear models into nonlinear models
- Recall that k(x,z) represents a dot product in some high-dim feature space ${\cal F}$
- Important: Any ML model/algo in which, during training and test, inputs only appear as dot product (pairwise similarity) can be "kernelized
- Just replace each term of the form $\mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j$ by $\phi(\mathbf{x}_i)^{\mathsf{T}}\phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) = K_{ij}$
- Most ML models/algos can be kernelized
- Let's look at an example: Kernelized SVM
 - Perhaps the most popular/natural example of kernelization



Kernelizing a Euclidean Distance

- Not just dot products but Eucliean distance can be kernelized too
- Many algorithms, e.g., LwP, KNN, etc. use Euclidean distances, e.g.,

$$d(a,b) = ||a - b||^2 = ||a||^2 + ||b||^2 - 2a^{\mathsf{T}}b = a^{\mathsf{T}}a + b^{\mathsf{T}}b - 2a^{\mathsf{T}}b$$

- This can be kernelized as well by replacing the above norms and inner products by their kernelized versions, assuming a kernel k with feature map ϕ

$$d(\phi(a), \phi(b)) = \|\phi(a) - \phi(b)\|^{2}$$

= $\phi(a)^{\top} \phi(a) + \phi(b)^{\top} \phi(b) - 2\phi(a)^{\top} \phi(b)$
= $k(a, a) + k(b, b) - 2k(a, b)$



Nonlinear SVM using Kernels



15

Kernelized SVM Training

Recall the soft-margin linear SVM objective (with no bias term)

- To kernelize, we can simply replace $G_{ij} = y_i y_j \mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j$ by $y_i y_j K_{ij}$
 - .. where $K_{ij} = k(x_i, x_j) = \phi(x_i)^{\mathsf{T}} \phi(x_j)$ for a suitable kernel function k
- The problem can now be solved just like the linear SVM case
- The new SVM learns a linear separator in kernel-induced feature space ${\cal F}$
 - This corresponds to a non-linear separator in the original feature space ${\mathcal X}$





Inputs only appear

as dot products \bigcirc

Kernelized SVM Prediction

- SVM weight vector for the kernelized case will be $w = \sum_{n=1}^{N} \alpha_n y_n \phi(x_n)$
- Imp: We can't store w unless the feature mapping $\phi(x_n)$ is finite dimensional
 - In practice, we store the α_n 's and the training data for test time (just like KNN)
 - In fact, need to store only training examples for which α_n is nonzero (i.e., the support vectors)
- Prediction for a new test input x_* (assuming hyperplane's bias b = 0) will be

$$y_* = \operatorname{sign}(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{\phi}(\boldsymbol{x}_*)) = \operatorname{sign}\left(\sum_{n=1}^{N} \alpha_n y_n \boldsymbol{\phi}(\boldsymbol{x}_n)^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x}_*)\right) = \operatorname{sign}\left(\sum_{n=1}^{N} \alpha_n y_n \boldsymbol{k}(\boldsymbol{x}_n, \boldsymbol{x}_*)\right)$$

- Note that the prediction cost also scales linearly with N (unlike a linear model where we only need to compute $w^T x_*$, whose cost only depends on D, not N)
- Also note that, for unkernelized (i.e., linear) SVM, $w = \sum_{n=1}^{N} \alpha_n y_n x_n$ can be computed and stored as a $D \times 1$ vector and we can compute $w^T x_*$ in O(D) time

Kernel extensions of other ML models



18

Kernel extensions of other ML models

- Most of the models what have studied can be kernelized
 - Kernel based linear/ridge regression
 - Kernel based LwP
 - Kernel based nearest neighbors
 - Kernel logistic regression
 - Kernel Perceptron

But the extra price has to be paid in terms of storage cost and slower predictions Kernel extension makes these approaches more powerful (nonlinear patterns can be learned)



CS771: Intro to ML

- Some of these extensions are simple to obtain, some not so (but possible)
- $\$ Imp: In these models, just like kernel SVM, the model parameters (e.g., the weight vector) can't be stored as a finite-dim vector (unless ϕ is finite dim)
 - Thus the training inputs need to be stored at test time as well

Also, just like kernel SVM, all of these will in general be slower at test time

Speeding-up Kernel Methods



20

Speeding-up Kernel Methods

Kernels assume that

$$k(\boldsymbol{x}_n, \boldsymbol{x}_m) = \phi(\boldsymbol{x}_n)^{\mathsf{T}} \phi(\boldsymbol{x}_m)$$

• Suppose for this kernel, we can get an L-dim feature vector $\psi(x)$ such that

$$k(\boldsymbol{x}_n, \boldsymbol{x}_m) \approx \boldsymbol{\psi}(\boldsymbol{x}_n)^{\mathsf{T}} \boldsymbol{\psi}(\boldsymbol{x}_m)$$

- Using features $\psi(x)$, we can learn a <u>linear model</u> with weights $w \in \mathbb{R}^{L}$
- This model will be a good approximation to the kernelized model
- Training will be faster because no need to store and work with kernel matrices
- Prediction at test time will also be faster we just need to compute $w^{\mathsf{T}}\psi(x_*)$
- Many ways to get such features $\psi(x)$ for standard kernels

Extracting Features using Kernels: Landmarks

• Suppose we choose a small set of L "landmark" inputs $z_1, z_2, ..., z_L$ in the training data



77

CS771: Intro to ML

 $\psi(\boldsymbol{x}_n) = [k(\boldsymbol{z}_1, \boldsymbol{x}_n), k(\boldsymbol{z}_2, \boldsymbol{x}_n), k(\boldsymbol{z}_3, \boldsymbol{x}_n)] \in \mathbb{R}^3$

• For each input x_n , using a kernel k, define an L-dimensional feature vector as follows

$$\psi(\boldsymbol{x}_n) = [k(\boldsymbol{z}_1, \boldsymbol{x}_n), k(\boldsymbol{z}_2, \boldsymbol{x}_n), \dots, k(\boldsymbol{z}_L, \boldsymbol{x}_n)] \in \mathbb{R}^L$$

- Can now apply a linear model on ψ representation (L-dimensional now) of the inputs
- This will be fast both at training as well as test time if L is small
- \blacksquare No need to kernelize the linear model while still reaping the benefits of kernels O

Extracting Feat. using Kernels: Random Features

Many kernel functions* can be written as

$$k(\boldsymbol{x}_n, \boldsymbol{x}_m) = \phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{x}_m) = \mathbb{E}_{\boldsymbol{w} \sim p(\boldsymbol{w})}[t_{\boldsymbol{w}}(\boldsymbol{x}_n)t_{\boldsymbol{w}}(\boldsymbol{x}_m)]$$

... where $t_w(.)$ is a function with params $w \in \mathbb{R}^D$ with w drawn from some distr. p(w)

- Example: For the RBF kernel, $t_w(.)$ is cosine func. and p(w) is zero mean Gaussian $k(x_n, x_m) = \mathbb{E}_{w \sim p(w)} [\cos(w^\top x_n) \cos(w^\top x_m)]$
- Given $w_1, w_2, ..., w_L$ from p(w), using Monte-Carlo approx. of above expectation $k(x_n, x_m) \approx \frac{1}{L} \sum_{\ell=1}^{L} \cos(w_{\ell}^{\top} x_n) \cos(w_{\ell}^{\top} x_m) = \psi(x_n)^{\top} \psi(x_m)$... where $\psi(x_n) = \frac{1}{\sqrt{L}} [\cos(w_1^{\top} x_n), ..., \cos(w_L^{\top} x_n)]$ is an *L*-dim vector
- Can apply a linear model on this L-dim rep. of the inputs (no need to kernelize)



Learning with Kernels: Some Aspects

- Storage/computational efficiency can be a bottleneck when using kernels
- During training, need to compute and store the $N \times N$ kernel matrix K in memory
- Need to store training data (or at least support vectors in case of SVMs) at test time
- Test time can be slow: O(N) cost to compute a quantity like $\sum_{n=1}^{N} \alpha_n k(\mathbf{x}_n, \mathbf{x}_*)$
- Approaches like landmark and random features can be used to speed up
- Choice of the right kernel is also very important
- Some kernels (e.g., RBF) work well for many problems but hyperparameters of the kernel function may need to be tuned via cross-validation
 Also, a lot of recent work on connections
- Quite a bit of research on learning the right kernel from data
 - Learning a combination of multiple kernels (Multiple Kernel Learning) learning
 - Bayesian kernel methods (e.g., Gaussian Processes) can learn the kernel hyperparameters from data(thus can be seen as learning the kernel)

between kernel

methods and deep

Deep Learning can also be seen as learning the kernel from data (more on this later) cs771: Intro to ML