



SOICT

HUST
ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

Kernelizing ML algorithms

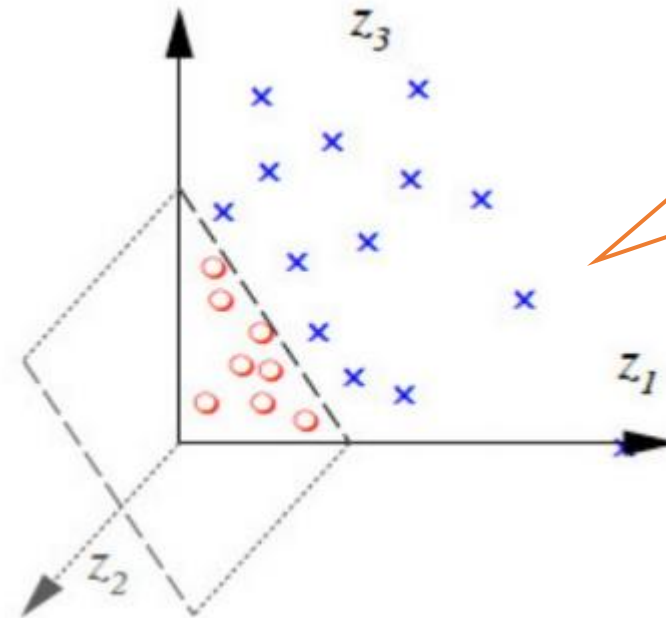
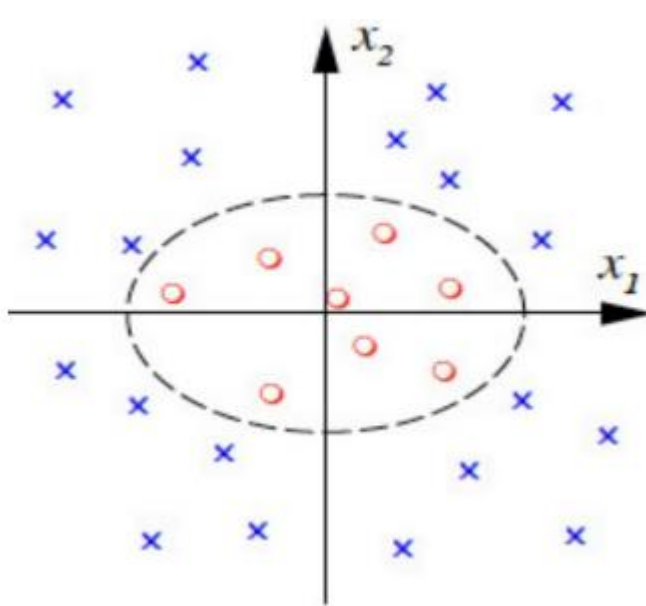
Dam Quang Tuan

Recap: Kernel functions

- Can assume a feature mapping ϕ that maps/transforms the inputs to a “nice” space

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



The linear model in the new feature space corresponds to a nonlinear model in the original feature space

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

Using Kernels

- Kernels can turn many linear models into nonlinear models
- Recall that $k(\mathbf{x}, \mathbf{z})$ represents a dot product in some high-dim feature space \mathcal{F}
- **Important:** Any ML model/algo in which, during training and test, inputs only appear as dot product can be “kernelized”
- Just replace each term of the form $\mathbf{x}_i^\top \mathbf{x}_j$ by $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) = K_{ij}$
- Most ML models/algos can be easily kernelized, e.g.,
 - Distance based methods, Perceptron, SVM, linear regression, etc.
 - Many of the unsupervised learning algorithms too can be kernelized (e.g., K-means clustering, Principal Component Analysis, etc. - will see later)
 - Let’s look at two examples: Kernelized SVM and Kernelized Ridge Regression

Nonlinear SVM using Kernels

Solving Soft-Margin SVM

- Recall the soft-margin SVM optimization problem

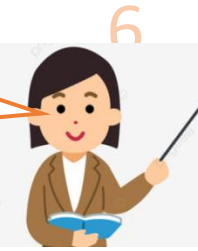
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & 1 \leq y_n(\mathbf{w}^T \mathbf{x}_n + b) + \xi_n, \quad -\xi_n \leq 0 \quad n = 1, \dots, N \end{aligned}$$

- Here $\xi = [\xi_1, \xi_2, \dots, \xi_N]$ is the vector of **slack variables**
- Introduce Lagrange multipliers α_n, β_n for each constraint and solve Lagrangian

$$\min_{\mathbf{w}, b, \xi} \max_{\alpha \geq 0, \beta \geq 0} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) - \xi_n\} - \sum_{n=1}^N \beta_n \xi_n$$

- The terms in red color above were not present in the hard-margin SVM
- Two set of dual variables $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_N]$ and $\beta = [\beta_1, \beta_2, \dots, \beta_N]$
- Will eliminate the primal var \mathbf{w}, b, ξ to get dual problem containing the dual variables

Solving Soft-Margin SVM



Note: if we ignore the bias term b then we don't need to handle the constraint $\sum_{n=1}^N \alpha_n y_n = 0$ (problem becomes a bit more easy to solve)

- The Lagrangian problem to solve

Otherwise, the α_n 's are coupled and some opt. techniques such as co-ordinate aspect can't easily applied

$$\min_{\mathbf{w}, b, \xi} \max_{\alpha \geq 0, \beta \geq 0} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) - \xi_n\} - \sum_{n=1}^N \beta_n \xi_n$$

- Take (partial) derivatives of \mathcal{L} w.r.t. \mathbf{w} , b , and ξ_n and setting to zero gives

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n, \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0, \quad \frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \Rightarrow C - \alpha_n - \beta_n = 0$$

Weighted sum of training inputs

- Using $C - \alpha_n - \beta_n = 0$ and $\beta_n \geq 0$, we have $\alpha_n \leq C$ (for hard-margin, $\alpha_n \geq 0$)

- Substituting these in the Lagrangian \mathcal{L} gives the Dual problem

The dual variables β don't appear in the dual problem!

Given α , \mathbf{w} and b can be found just like the hard-margin SVM case

$$\max_{\alpha \leq C, \beta \geq 0} \mathcal{L}_D(\alpha, \beta) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n y_m y_n (\mathbf{x}_m^T \mathbf{x}_n) \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

Maximizing a concave function (or minimizing a convex function) s.t. $\alpha \leq C$ and $\sum_{n=1}^N \alpha_n y_n = 0$. Many methods to solve it

$$\max_{\alpha \leq C} \mathcal{L}_D(\alpha) = \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{G} \alpha$$

In the solution, α will still be sparse just like the hard-margin SVM case. Nonzero α_n correspond to the support vectors

(Note: For various SVM solvers, can see "Support Vector Machine Solvers" by Bottou and Lin)

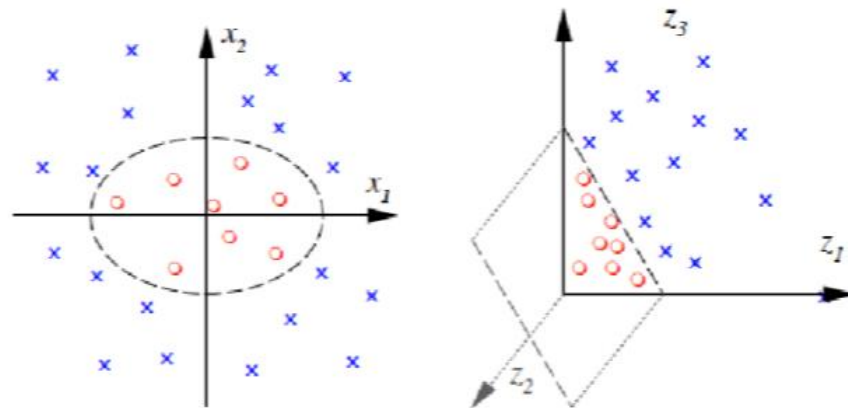
Kernelized SVM Training

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

$$\operatorname{argmax}_{\mathbf{0} \leq \boldsymbol{\alpha} \leq \mathbf{C}} \quad \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{G} \boldsymbol{\alpha}$$

$G_{ij} = y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$

- To kernelize, we can simply replace $G_{ij} = y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$ by $y_i y_j K_{ij}$
 - .. where $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ 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 \mathcal{F}
 - This corresponds to a **non-linear separator** in the original feature space \mathcal{X}



Kernelized SVM Prediction

- SVM weight vector for the kernelized case will be $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \phi(\mathbf{x}_n)$
- Note: We can't store \mathbf{w} unless the feature mapping $\phi(\mathbf{x}_n)$ 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 \mathbf{x}_* (assuming hyperplane's bias $b = 0$) will be

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

Nonlinear Ridge Regression using Kernels

Kernelized Ridge Regression

- Recall the ridge regression problem $\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \mathbf{w}^\top \mathbf{w}$

- The solution to this problem was

Inputs don't appear to be as inner product. No hope of kernelization?



They do; with a bit of algebra ☺

$$\mathbf{w} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right) \left(\sum_{n=1}^N y_n \mathbf{x}_n \right) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Can use matrix inversion lemma $(\mathbf{F}\mathbf{H}^{-1}\mathbf{G} - \mathbf{E})^{-1}\mathbf{F}\mathbf{H}^{-1} = \mathbf{E}^{-1}\mathbf{F}(\mathbf{G}\mathbf{E}^{-1}\mathbf{F} - \mathbf{H})^{-1}$

- Using the lemma, can rewrite \mathbf{w} as

$$\mathbf{w} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \mathbf{X}^\top \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \mathbf{x}_n \quad \text{where} \quad \boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

$N \times 1$ vector of dual variables

Note: Not sparse unlike SVM

- Kernelized weight vector will be $\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$

Prediction cost is also linear in N (like KNN)

- Prediction for a test input \mathbf{x}_* will be $\mathbf{w}^\top \phi(\mathbf{x}_*) = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_*) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}_*)$

Speeding-up Kernel Methods

Speeding-up Kernel Methods

- Kernel methods, unlike linear models are slow at training and test time
- Would be nice if we could easily compute mapping $\phi(\mathbf{x})$ associated with kernel k
- Then we could apply linear models directly on $\phi(\mathbf{x})$ without having to kernelize
- But this is in general not possible since $\phi(\mathbf{x})$ is very high/infinite dimensional
- An alternative: Get a good set of low-dim features $\psi(\mathbf{x}) \in \mathbb{R}^L$ using the kernel k
- If $\psi(\mathbf{x})$ is a good approximation to $\phi(\mathbf{x})$ then we can use $\psi(\mathbf{x})$ in a linear model

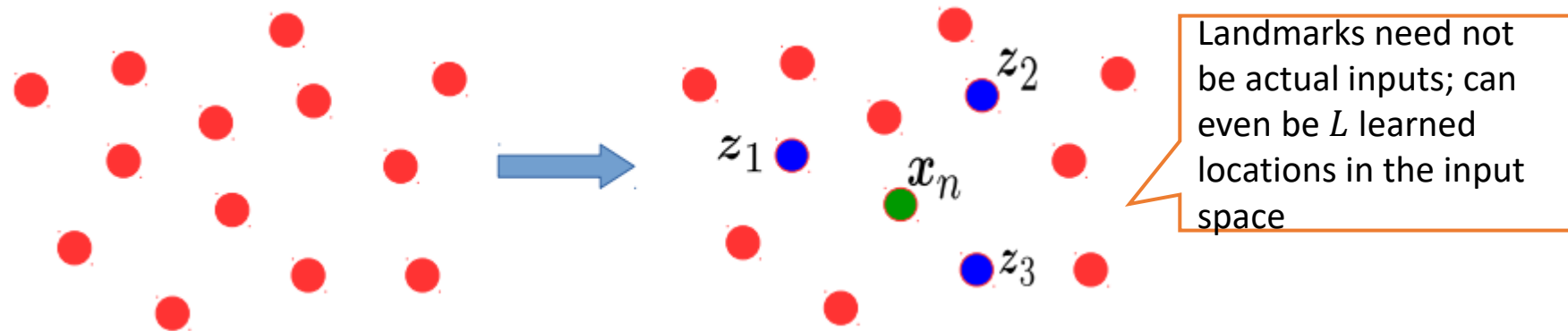
Goodness Criterion:
$$\psi(\mathbf{x}_i)^\top \psi(\mathbf{x}_j) \approx \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

... which also means $\psi(\mathbf{x}_i)^\top \psi(\mathbf{x}_j) \approx k(\mathbf{x}_i, \mathbf{x}_j)$

- Will look at two popular approaches: [Landmarks](#) and [Random Features](#)

Extracting Features using Kernels: Landmarks

- Suppose we choose a small set of L “landmark” inputs $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_L$ in the training data



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

- For each input \mathbf{x}_n , using a kernel k , define an L -dimensional feature vector as follows

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

- 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

Extracting Features using Kernels: Random Features

- Many kernel functions* can be written as

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

.. where $t_{\mathbf{w}}(\cdot)$ is a function with params $\mathbf{w} \in \mathbb{R}^L$ with \mathbf{w} drawn from some distr. $p(\mathbf{w})$

- Example: For the RBF kernel, $t_{\mathbf{w}}(\cdot)$ is cosine func. and $p(\mathbf{w})$ is zero mean Gaussian

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

- Given $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_L$ from $p(\mathbf{w})$, using Monte-Carlo approx. of above expectation

$$k(\mathbf{x}_n, \mathbf{x}_m) \approx \frac{1}{L} \sum_{\ell=1}^L \cos(\mathbf{w}_\ell^\top \mathbf{x}_n) \cos(\mathbf{w}_\ell^\top \mathbf{x}_m) = \psi(\mathbf{x}_n)^\top \psi(\mathbf{x}_m)$$

.. where $\psi(\mathbf{x}_n) = \frac{1}{\sqrt{L}} [\cos(\mathbf{w}_1^\top \mathbf{x}_n), \dots, \cos(\mathbf{w}_L^\top \mathbf{x}_n)]$ 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 \mathbf{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
- Quite a bit of research on learning the right kernel from data
 - Learning a combination of multiple kernels ([Multiple Kernel Learning](#))
 - [Bayesian kernel methods](#) (e.g., [Gaussian Processes](#)) can learn the kernel hyperparameters from data (thus can be seen as learning the kernel)
 - Deep Learning can also be seen as learning the kernel from data (more on this later)

References

CS771: Intro to Machine Learning (Fall 2021), Nisheeth Srivastava, IIT Kanpur