

Kernel Methods

Sargur Srihari
srihari@buffalo.edu

Topics in Kernel Methods

1. Linear Models vs Memory-based models
2. Stored Sample Methods
3. Kernel Functions
 - Dual Representations
 - Constructing Kernels
4. Extension to Symbolic Inputs
5. Fisher Kernel

Linear Models vs Memory-based models

- Linear parametric models for regression and classification have the form $y(\mathbf{x}, \mathbf{w})$
 - During learning we either get a maximum likelihood estimate of \mathbf{w} or a posterior distribution of \mathbf{w}
 - Training data is then discarded
 - Prediction based only on vector \mathbf{w}
 - This is true of Neural networks as well
- Memory-based methods keep the training samples (or a subset) and use them during the prediction phase

Memory-Based Methods

- Training data points are used in prediction
- Examples of such methods
 - Parzen probability density model
 - Linear combination of kernel functions centered on each training data point
 - Nearest neighbor classification
- These are memory-based methods
- They require a metric to be defined
 - Fast to train, slow to predict

Kernel Functions

- Linear models can be re-cast
 - into equivalent dual where predictions are based on kernel functions evaluated at training points

- Kernel function is given by

$$k(x, x') = \phi(x)^T \phi(x')$$

- where $\phi(x)$ is a fixed nonlinear mapping (basis function)
- Kernel is a symmetric function of its arguments

$$k(x, x') = k(x', x)$$

- Kernel can be interpreted as similarity of x and x'
 - Simplest is identity mapping in feature space $\phi(x) = x$
 - In which case $k(x, x') = x^T x'$
 - Called Linear Kernel

Kernel Trick

- Formulated as inner product allows extending well-known algorithms
 - by using the kernel trick
- Basic idea of kernel trick
 - If an input vector x appears only in the form of scalar products then we can replace scalar products with some other choice of kernel
- Used widely
 - in support vector machines
 - in developing non-linear variant of PCA
 - In kernel Fisher discriminant

Other Forms of Kernel Functions

- Function of difference between arguments

$$k(x, x') = k(x - x')$$

- Called *stationary* kernel since invariant to translation

- Radial basis functions

$$k(x, x') = k(\|x - x'\|)$$

- Depends on magnitude of distance between arguments
 - Note that the kernel function is scalar while x is M -dimensional

For these to be valid kernel functions they should be shown to have the property

$$k(x, x') = \varphi(x)^T \varphi(x')$$

Dual Representation

- Linear models for regression/classification can be reformulated in terms of dual representation
 - In which kernel function arises naturally
 - Plays important role in SVMs
- Consider linear regression model
 - Parameters from minimizing *sum-of-squares*

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

where $\mathbf{w} = (w_0, \dots, w_{M-1})^T$, $\phi = (\phi_0, \dots, \phi_{M-1})^T$
we have N samples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
 λ is the regularization coefficient

ϕ is the set of M
basis functions
or feature vector

- Minimum obtained by setting gradient of $J(\mathbf{w})$ wrt \mathbf{w} equal to zero

Solution for \mathbf{w} : linear combination of $\phi(\mathbf{x}_n)$

- Equating derivative $J(\mathbf{w})$ wrt \mathbf{w} to zero we get

$$\begin{aligned}\mathbf{w} &= -\frac{1}{\lambda} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \phi(\mathbf{x}_n) \\ &= \sum_{n=1}^N a_n \phi(\mathbf{x}_n) \\ &= \Phi^T \mathbf{a}\end{aligned}$$

- Solution for \mathbf{w} : linear combination of $\phi(\mathbf{x}_n)$ whose coefficients are functions of \mathbf{w} where
 - Φ is design matrix whose n^{th} row is given by $\phi(\mathbf{x}_n)^T$

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \cdot & \cdot & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \cdot & \cdot & \phi_{M-1}(x_2) \\ \vdots & & & \vdots \\ \phi_0(x_N) & \cdot & \cdot & \phi_{M-1}(x_N) \end{bmatrix} \text{ is a } N \times M \text{ matrix}$$

- Vector $\mathbf{a} = (a_1, \dots, a_N)^T$ has the definition

$$a_n = -\frac{1}{\lambda} \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}$$

Transformation from w to a

- Thus we have $w = \Phi^T a$
- Instead of working with parameter vector w we can reformulate least squares algorithm in terms of parameter vector a
 - giving rise to dual representation

$$a_n = -\frac{1}{\lambda} \left\{ w^T \phi(x_n) - t_n \right\}$$

- Although the definition of a still includes w
It can be eliminated by the use of the kernel function

Gram Matrix and Kernel Function

- Gram matrix $K = \Phi \Phi^T$ is $N \times N$ matrix
 - with elements

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

- where kernel function $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdot & \cdot & k(\mathbf{x}_1, \mathbf{x}_N) \\ \cdot & & & \\ \cdot & & & \\ k(\mathbf{x}_N, \mathbf{x}_1) & & & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

Gram Matrix Definition:
Given N vectors, it is the matrix of all inner products

- Notes:**
 - Φ is $N \times M$ and K is $N \times N$ **Note:** $N \times M$ times $M \times N$
 - K is a matrix of similarities of pairs of samples (thus it is symmetric)

Error Function in Terms of Gram Matrix

- Error Function is

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

- Substituting $\mathbf{w} = \Phi^T \mathbf{a}$ into $J(\mathbf{w})$ gives

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

where $\mathbf{t} = (t_1, \dots, t_N)^T$

- Error function is written in terms of Gram matrix as

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

- Solving for \mathbf{a} by combining $\mathbf{w} = \Phi^T \mathbf{a}$ and $a_n = -\frac{1}{\lambda} \left\{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \right\}$

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

Solution for \mathbf{a} can be expressed as a linear combination of elements of $\phi(\mathbf{x})$ whose coefficients are entirely in terms of kernel $k(\mathbf{x}, \mathbf{x}')$ from which we can recover original formulation in terms of parameters \mathbf{w}

Prediction Function

- Prediction for new input \mathbf{x}
 - We can write $\mathbf{a} = (K + \lambda I_N)^{-1} \mathbf{t}$ by combining $\mathbf{w} = \Phi^T \mathbf{a}$ and

$$a_n = -\frac{1}{\lambda} \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}$$

- Substituting back into linear regression model,

$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) \\ &= \mathbf{a}^T \Phi \phi(\mathbf{x}) \\ &= k(\mathbf{x})^T (K + \lambda I_N)^{-1} \mathbf{t} \quad \text{where } k(\mathbf{x}) \text{ has elements } k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x}) \end{aligned}$$

- Prediction is a linear combination of the target values from the training set.

Advantage of Dual Representation

- Solution for a is entirely in terms of $k(x, x')$
 - From a we can recover w using $w = \Phi^t a$
- In parametric formulation, solution is $w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$
 - We invert $M \times M$ matrix, since Φ is $N \times M$
- In dual formulation, solution is $a = (K + \lambda I_N)^{-1} t$
 - We are inverting an $N \times N$ matrix
 - An apparent disadvantage, since $N \gg M$
 - Advantage of dual: work with kernel $k(x, x')$, so:
 - avoid working with a feature vector $\phi(x)$ and
 - problems associated with very high or infinite dimensionality of x

Constructing Kernels

- To exploit kernel substitution need valid kernels
 - Two methods:
 - (1) from $\phi(x)$ to $k(x, x')$
 - (2) from $k(x, x')$ to $\phi(x)$
 - First Method
 - Choose $\phi(x)$ and use it to find corresponding kernel

$$\begin{aligned}k(x, x') &= \phi(x)^T \phi(x') \\ &= \sum_{i=1}^M \phi_i(x) \phi_i(x')\end{aligned}$$

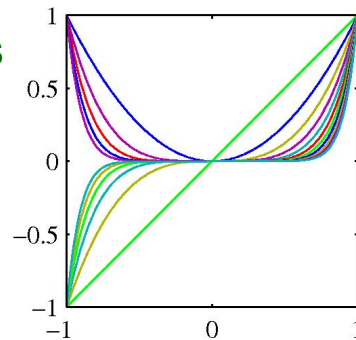
- where $\phi_i(x)$ are basis functions
 - One-dimensional input space is shown next

Kernel Functions from basis functions

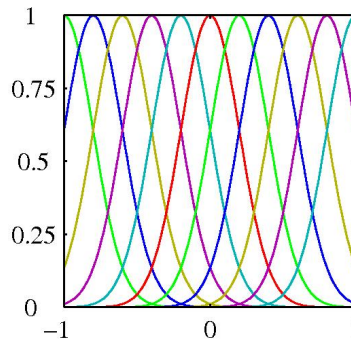
One-dimensional input space

Basis
Functions
 $\phi_i(x)$

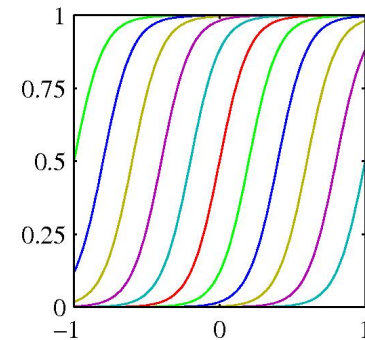
Polynomials



Gaussian

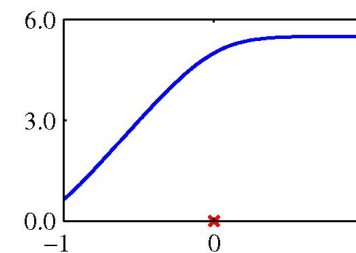
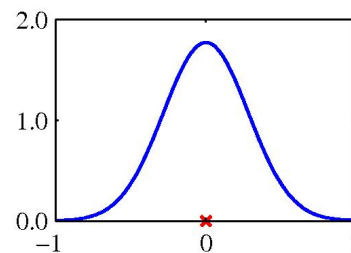
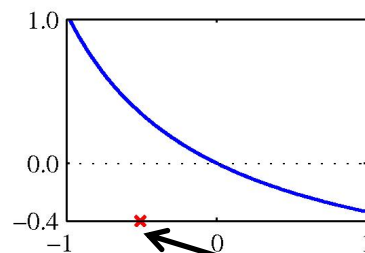


Logistic Sigmoid



Kernel
Functions

$$k(x, x') = \phi(x)^T \phi(x')$$



Red cross is x'

Method 2: Direct Construction of Kernels

- Kernel has to correspond to a scalar product in some (perhaps infinite dimensional) space

- Consider kernel function $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$

- In 2-D space

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

- Feature mapping takes the form $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)$
 - Comprises of all second order terms with a specific weighting
 - Inner product needs computing 6 feature values, 9 multiplications
- Kernel function $k(\mathbf{x}, \mathbf{z})$ has 2 multiplies and a squaring
- For $(\mathbf{x}^T \mathbf{z} + c)^2$ we get constant, linear, second order terms
- For $(\mathbf{x}^T \mathbf{z} + c)^M$ we get all powers of \mathbf{x} (monomials)

Testing whether a function is a valid kernel

- Without constructing $\phi(\mathbf{x})$, Necessary and sufficient condition for $k(\mathbf{x}, \mathbf{x}')$ to be a kernel is
 - Gram matrix K , with elements $k(\mathbf{x}_n, \mathbf{x}_m)$ is positive $\{\mathbf{x}_n\}$ semi-definite for all possible choices of the set
 - Positive semi-definite is not the same thing as a matrix whose elements are non-negative
 - It means $\mathbf{z}^T K \mathbf{z} \geq 0$ for non-zero vectors \mathbf{z} with real entries
 i.e., $\sum_n \sum_m K_{nm} z_n z_m \geq 0$ for any real numbers z_n, z_m
 - Mercer's theorem: any continuous, symmetric, positive semi-definite kernel function $k(\mathbf{x}, \mathbf{y})$ can be expressed as a dot product in a high-dimensional space
- New kernels can be constructed from simpler kernels as building blocks

Techniques for Constructing Kernels

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

- $k(\mathbf{x}, \mathbf{x}') = c k_1(\mathbf{x}, \mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) k_1(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$
- $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$
- $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = k_3(f(\mathbf{x}), f(\mathbf{x}'))$
- $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}'$
- $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}_a') + k_b(\mathbf{x}_b, \mathbf{x}_b')$
- $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}_a') k_b(\mathbf{x}_b, \mathbf{x}_b')$

Where

$f(\cdot)$ is any function

$q(\cdot)$ is a polynomial with non-negative coefficients

$f(\mathbf{x})$ is a function from \mathbf{x} to \mathbb{R}^M

k_3 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 with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$

k_a and k_b are valid kernel functions

Kernels for specific applications

- Requirements for $k(\mathbf{x}, \mathbf{x}')$
 - It is symmetric
 - Its Gram matrix is positive semidefinite
 - It expresses the appropriate similarity between \mathbf{x} and \mathbf{x}' for the intended application

Gaussian Kernel

- Commonly used kernel is

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

- It is seen as a valid kernel by expanding square

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

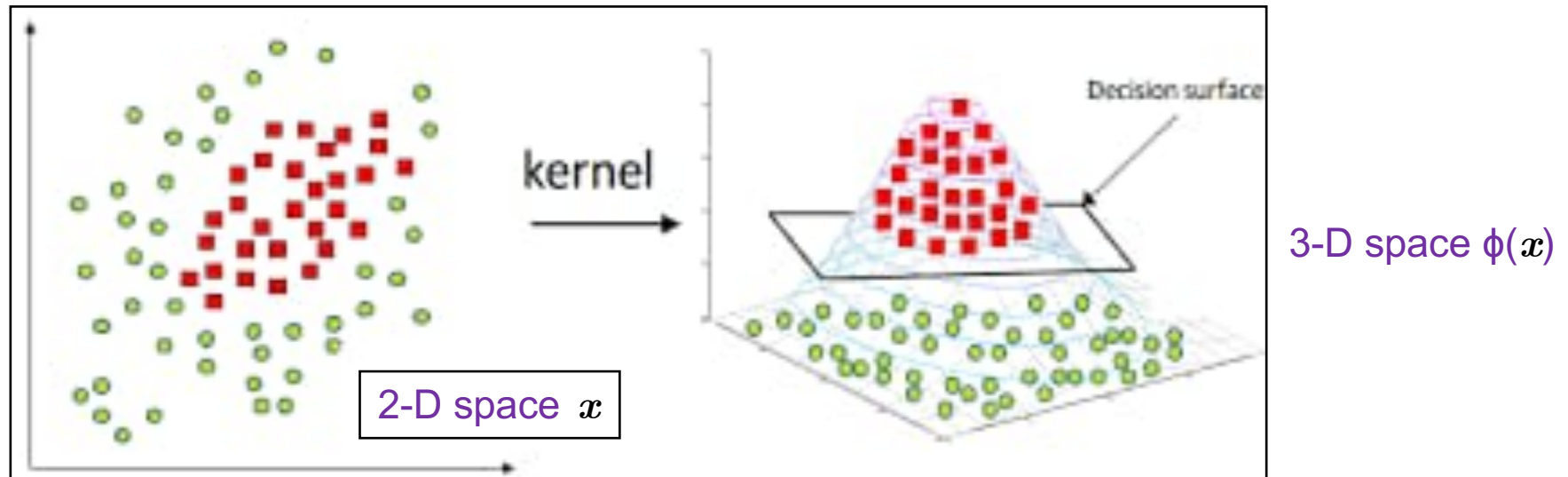
- To give

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

- From kernel construction rules 2 and 4
 - together with validity of linear kernel $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$
- Can be extended to non-Euclidean distances

$$k(\mathbf{x}, \mathbf{x}') = \exp \{ (-1/2\sigma^2) [\kappa(\mathbf{x}, \mathbf{x}') + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}')] \}$$

Use of a kernel method



→ Infinite dimensional space $\phi(x)$
Induced by a Gaussian kernel

Oil slick 2D- image
that looks like a feature space x

Extension of Kernels to Symbolic Inputs

- Important contribution of kernel viewpoint:
 - Inputs are symbolic, not vectors of real numbers
 - Kernels defined for graphs, sets, strings, text documents
- If A_1 and A_2 are two subsets of objects
 - A simple kernel is $k(A_1, A_2) = 2^{|A_1 \cap A_2|}$
 - where $|A|$ indicates no of elements in A
 - A valid kernel since it can be shown to correspond to an inner product in a feature space

$$A = \{1, 2, 3, 4, 5\}$$

$$A_1 = \{2, 3, 4, 5\}, A_2 = \{1, 2, 4, 5\}, A_1 \cap A_2 = \{2, 4, 5\}$$

$$\text{Hence } k(A_1, A_2) = 8$$

Example feature vectors:

$$\phi(A_1) = [2 \ 2 \ 2 \ 2] \text{ and } \phi(A_2) = [1 \ 1 \ 1 \ 1]$$

such that

$$\phi(A_1)^T \phi(A_2) = 8$$

Kernels for Complex Objects

- ML methods studied so far require input represented as a fixed-size feature vector $x_i \in \mathbb{R}^D$
 - For certain objects it is not clear how best to represent them as feature vectors, E.g.,
 - Text or protein sequence of variable length
 - Molecular structure with complex 3-D geometry
 - Evolutionary tree which has variable size and shape
- Solutions
 - Define generative model, whose parameters are features. Plug these features into standard models
 - Measure similarity between objects which does not require preprocessing into feature vector format

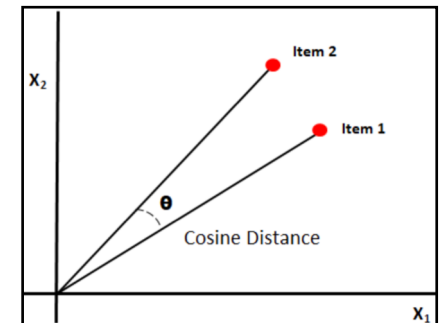
Kernels for Comparing Documents

- In IR and document classification, need a way of comparing documents \mathbf{x}_i and $\mathbf{x}_{i'}$,
- Bag of words representation: x_{ij} is the no of times word j occurs in document i

$$\mathbf{x}_i = [x_{i1}, \dots, x_{iD}]$$

- Cosine similarity

$$k(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\mathbf{x}_i^T \mathbf{x}_{i'}}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_{i'}\|_2}$$



- It measures cosine of angle between \mathbf{x}_i and $\mathbf{x}_{i'}$,
- Since \mathbf{x}_i is a count vector (and hence non-negative) the cosine similarity is between 0 and 1
 - Where 0 means the vectors are orthogonal and therefore have no words in common

Disadvantage of Cosine Similarity

- It does not work well for two reasons
 1. Stop words
 - If x_i has any word in common with x_j , it is deemed similar, even though some words such as the, and occur commonly in many documents and are not discriminative
 2. High frequency words in a document
 - If a discriminatory word occurs many times in a document, the similarity is artificially boosted
 - Even though word usage would be bursty, i.e., once a word is used in a document it is very likely to be used again

TF-IDF

- Cosine similarity performance can be improved with some preprocessing
- Use feature vector called TF-IDF representation
 - Term frequency-Inverse document frequency

- TF is a log-transform of the count

$$\text{tf}(x_{ij}) \triangleq \log(1 + x_{ij})$$

- IDF is defined as

$$\text{idf}(j) \triangleq \log \frac{N}{1 + \sum_{i=1}^N I(x_{ij} > 0)}$$

x_{ij} is the no of times word j occurs in document i

- Where N is the total no of documents
- Denominator counts how many documents contain term j

- Finally we define

$$\text{tf-idf}(xi) \triangleq [\text{tf}(x_{ij}) \times \text{idf}(j)]_{j=1}^V$$

- We use this inside cosine similarity measure

$$k(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_{i'})}{\|\phi(\mathbf{x}_i)\|_2 \|\phi(\mathbf{x}_{i'})\|_2}$$

- Where $\phi(\mathbf{x}) = \text{tf-idf}(\mathbf{x})$

String Kernels

- Real power of kernels: inputs are structured
- Compare two variable-length strings
 - Strings x , and x' of length D , D' defined over A
 - E.g., two amino acid sequences defined over 20-letter alphabet $A=\{A,R,N,D,C,E,Q,G,H,I,L,K,M,F,P,S,T,W,Y,V\}$
 - x is the sequence of length 110
 IPTSALVKETLALLSTHRTLIIANETLRIPVPVHKN.....VNQFLDY**LQE**FLGVMNTEWI
 - x' is a string of length 153
 PHRRDLCSRSLWLARKIRSDLTALTESYVKHQGLWSELTEAER**LQEN**LQAYRTFHVLLA.....
 - The strings have the substring **LQE** in common
 - We can define the similarity of two strings to be the no of substrings they have in common
 - Definition given next

Mercer Kernel

- If s is a substring of x we can write $x=usv$ for some (possibly empty) strings u, s and v
- $\phi_s(x)$ is no of times substring s appears in x
- Define kernel between two strings x and x' as

$$k(x, x') = \sum_{s \in A^*} w_s \phi_s(x) \phi_s(x')$$

- Where $w_s \geq 0$ and A^* is the set of all strings from alphabet A (known as Kleene * operator)
- Can be computed in $O(|x| + |x'|)$ time
- Cases of interest:
 - $w_s = 0$ for $|s| > 1$ we get bag-of-characters kernel— defines $\phi(x)$ as no of times each character in A occurs in x
 - If we require s to be bordered by white space we get bag-of- words kernel

Combining Discriminative and Generative Models

- Generative models deal naturally with missing data and with HMM of varying length
- Discriminative models such as SVM have better performance
- Can use a generative model to define a kernel and use kernel in discriminative approach

Kernels based on Generative Models

- For a model $p(\mathbf{x})$ we define a kernel by

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

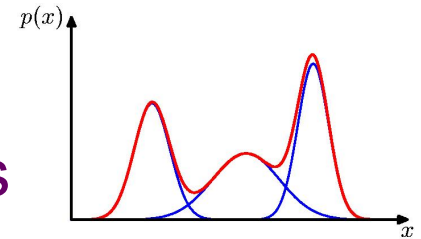
- A valid kernel since an inner product in 1-dimensional feature space defined by the mapping $p(\mathbf{x})$
- Two inputs \mathbf{x} and \mathbf{x}' are similar if they have high probabilities

Kernel Functions based on Mixture Densities

- Extension to sums/products of distributions

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

- where $p(i)$ are positive weighting coefficients
- A valid kernel based on two rules of kernel construction:



$$k(\mathbf{x}, \mathbf{x}') = c k_1(\mathbf{x}, \mathbf{x}') \text{ and } k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

- Two inputs \mathbf{x} and \mathbf{x}' will give a large value of k , and hence appear similar, if they have a significant probability under a range of different components
- Taking the limit to infinite sum

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

- where \mathbf{z} is a continuous latent variable

Kernels for Sequences

- Data consists of ordered sequences of length L

$$X = \{x_1, \dots, x_L\}$$

- Generative model for sequences is HMM
 - Hidden states $Z = \{z_1, \dots, z_L\}$
- Kernel Function for measuring similarity of sequences X and X' is

$$k(X, X') = \sum_Z p(X | Z) p(X' | Z) p(Z)$$

- Both observed sequences are generated by same hidden sequence Z

Fisher Kernel

- Alternative technique for generative models
 - In document retrieval, protein sequences
- Consider parametric generative model $p(\mathbf{x}|\boldsymbol{\theta})$ where $\boldsymbol{\theta}$ denotes vector of parameters
 - Goal: find kernel that measures similarity of two vectors \mathbf{x} and \mathbf{x}' induced by the generative model
- Define Fisher score as gradient wrt $\boldsymbol{\theta}$

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x} | \boldsymbol{\theta})$$

A vector of same dimensionality as $\boldsymbol{\theta}$

- **Fisher Kernel is** $k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^T \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}')$

Fisher score is more generally the gradient of the log-likelihood

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} [\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^T]$$

where \mathbf{F} is the Fisher information matrix

Fisher Information Matrix

- Presence of Fisher information matrix causes kernel to be invariant under non-linear parametrization of the density model $\theta \rightarrow \psi(\theta)$
 - In practice, infeasible to evaluate Fisher Information Matrix. Instead use the approximation

$$F \approx \frac{1}{N} \sum_{n=1}^N g(\theta, \mathbf{x}_n) g(\theta, \mathbf{x}_n)^T \quad k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T F^{-1} g(\theta, \mathbf{x}')$$

- This is the covariance matrix of the Fisher scores
 - So the Fisher kernel corresponds to whitening of the Fisher scores $k(\mathbf{x}, \mathbf{x}') = g(\theta, \mathbf{x})^T g(\theta, \mathbf{x}')$
- More simply omit F and use non-invariant kernel

Sigmoidal Kernel

- A link between SVMs and neural network

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\mathbf{a} \mathbf{x}^T \mathbf{x}' + b)$$

- Its Gram matrix is not positive semidefinite
- But used in practice because it gives SVMs a superficial resemblance to neural networks
- Bayesian neural network with an appropriate prior reduces to a Gaussian process
 - Provides a deeper link between neural networks and kernel methods