



Kernels

Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Andrew Ng (Stanford), Eric Eaton (UPenn), Tommi Jaakola (MIT) and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

SVMs

Learning Goals

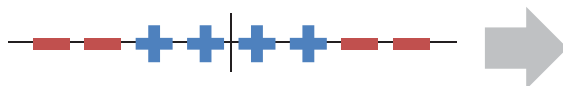
- ✓ Solving the SVM optimization problem
 - using Lagrange multipliers (leads to dual problem)
- ✓ Allowing misclassified examples
 - using slack variables (leads to soft-margin SVM)
- ✓ Describe the SVM loss function
- Allowing non-linear decision boundaries
 - using kernels

Kernel Basics

Learning Goals

- Motivate kernels
 - mapping to new feature space
 - easy to use in SVM optimization and prediction
- Define kernels formally
 - what makes a kernel valid

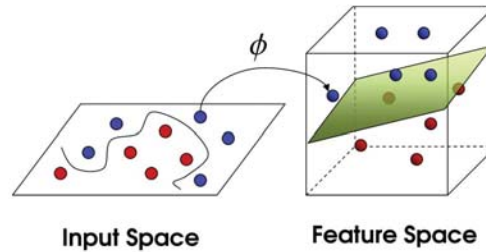
When Linear Separators Fail



$x \in \mathbb{R}$
not linearly separable



Mapping to a New Feature Space



$$\phi : \mathcal{X} \mapsto \mathcal{F}$$

Example: for $\mathbf{x} \in \mathbb{R}^2$,

$$\phi([x_1, x_2]^T) = [x_1, x_2, x_1 x_2, x_1^2, x_2^2]^T$$

Rather than run SVM on \mathbf{x} , run it on $\phi(\mathbf{x})$

Find non-linear separator in input space

Based on slide by Eric Eaton [originally by Tim Oates; image from <http://web.engr.oregonstate.edu/~afern/classes/cs534/>]

Using Mapped Features for SVM

Primal
$$\min_{\boldsymbol{\theta}, b} \frac{1}{2} \|\boldsymbol{\theta}\|^2$$

s.t.
$$\mathbf{y}^{(i)} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} + b) \geq 1 \text{ for } i = 1, \dots, n$$

Dual
$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle$$

s.t.
$$\alpha_i \geq 0 \quad \text{for } i = 1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0$$

Solution

$$\boldsymbol{\theta} = \sum_{i=1}^n \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

for SV $\mathbf{x}^{(i)}$,
$$b = \mathbf{y}^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}$$

Prediction

$$\hat{y} = \boldsymbol{\theta}^T \mathbf{x} + b$$

$$= \left(\sum_{i=1}^n \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \mathbf{x} + b$$

$$= \sum_{i=1}^n \alpha_i y^{(i)} \langle \mathbf{x}^{(i)}, \mathbf{x} \rangle + b$$

If we have found optimal α_i 's (from dual formulation), then to make a prediction for \mathbf{x} , we only have to calculate a quantity dependent on dot product between \mathbf{x} and SVs in training set.

Key Insight

Optimization

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle$$

s.t. $\alpha_i \geq 0$ for $i = 1, \dots, n$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0$$

optimization depends only on inner products between inputs $\langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle$

Prediction

$$\hat{y} = \sum_{i=1}^n \alpha_i y^{(i)} \langle \mathbf{x}^{(i)}, \mathbf{x} \rangle + b$$

prediction depends only on inner products between test example and SVs $\langle \mathbf{x}^{(i)}, \mathbf{x} \rangle$

Kernels

- Capture nonlinear patterns in data
 - because linear models (e.g. regression, SVM) may not be rich enough
- **kernels** make linear models work in non-linear settings
 - by mapping to higher dimensions $\mathbf{x} \rightarrow \phi(\mathbf{x})$
 - and applying linear model in new input space
- Problem
 - computing mapping may be inefficient
 - using mapped features could be inefficient
- Solution: kernels!
 - mapping does not have to be explicitly computed
 - computations with mapped features remain efficient

The Quadratic Kernel

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$



(This slide intentionally left blank.)

Formal Definition

$$\phi : \mathcal{X} \mapsto \mathcal{F}$$

$$k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$$

ϕ takes input $x \in \mathcal{X}$ (input space)
and maps to \mathcal{F} (feature space)

kernel k takes two inputs $x \in \mathcal{X}$ and $z \in \mathcal{X}$
and computes their similarity $\langle \phi(x), \phi(z) \rangle$ in \mathcal{F}

Can any function be used as a kernel function?

Note: \mathcal{F} needs to be a vector space with a dot product defined on it (aka a Hilbert space).



Valid Kernels

Kernel Matrix (aka Gram Matrix)

Kernel k also defines kernel matrix \mathbf{K} over data

Given m examples (m finite, not necessarily training set),
let square $m \times m$ matrix \mathbf{K} be defined so that

$$K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$$

Notes

- \mathbf{K} is a $m \times m$ matrix of pairwise similarities
- $K_{ij} = K_{ji}$ (by symmetry of dot products) so \mathbf{K} is symmetric

Theorem (Mercer)

Let $k : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$ be given. Then for k to be a valid kernel, it is necessary and sufficient that for any $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ (n finite), the corresponding kernel matrix \mathbf{K} is symmetric positive semi-definite. [prove in homework]

Reminder: A matrix \mathbf{K} is PSD if for all real vectors α , $\alpha^T \mathbf{K} \alpha \geq 0$.

Kernels

Learning Goals

- Describe common kernels
 - linear, polynomial, Gaussian (RBF)
- Prove that RBF kernel is a valid kernel
 - using kernel closure properties

Polynomial Kernel

Let $k(\mathbf{x}, \mathbf{z}) = (1 + \langle \mathbf{x}, \mathbf{z} \rangle)^p$ (hyperparameter $p = 1, 2, \dots$).

Then $\phi(\mathbf{x})$ contains all terms up to degree p .

Q: For $k(\mathbf{x}, \mathbf{z}) = (1 + \langle \mathbf{x}, \mathbf{z} \rangle)^2$, where $\mathbf{x}, \mathbf{z} \in \mathbb{R}^2$,
what is the corresponding mapping $\phi(\mathbf{x})$?

S:

scikit-learn: $k(\mathbf{x}, \mathbf{z}) = (\gamma \langle \mathbf{x}, \mathbf{z} \rangle + r)^d$

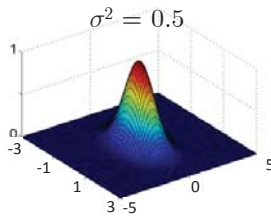
γ and r trade-off influence of lower-order terms



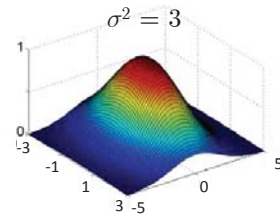
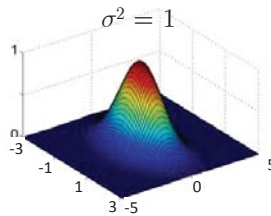
Gaussian Kernel (aka Radial Basis Function Kernel)

$$k(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{z}\|_2^2}{2\sigma^2}\right) \quad (\text{hyperparameter } \sigma^2 > 0)$$

- Has value 1 when $\mathbf{x} = \mathbf{z}$, value falls off to 0 with increasing distance
- Interpreting σ^2



less "smooth" decision boundary
 σ^2 too small \Rightarrow overfit



smoother decision boundary
 σ^2 too large \Rightarrow underfit

Note: need to do feature scaling before using Gaussian Kernel

Is this a valid kernel?

scikit-learn: $k(\mathbf{x}, \mathbf{z}) = \exp(-\gamma\|\mathbf{x}-\mathbf{z}\|^2)$, where $\gamma = 1/(2\sigma^2) > 0$

Images from Eric Eaton



Popular Kernels

- linear $k(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle$
- polynomial $k(\mathbf{x}, \mathbf{z}) = (\gamma\langle \mathbf{x}, \mathbf{z} \rangle + r)^d$
- Gaussian (RBF) $k(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{z}\|_2^2}{2\sigma^2}\right)$
- sigmoid $k(\mathbf{x}, \mathbf{z}) = \tanh(\gamma\langle \mathbf{x}, \mathbf{z} \rangle + r)$
 - SVM with sigmoid kernel equivalent to 2-layer perceptron (neural network)
- cosine $k(\mathbf{x}, \mathbf{z}) = \frac{\langle \mathbf{x}, \mathbf{z} \rangle}{\|\mathbf{x}\|\|\mathbf{z}\|} = \left\langle \frac{\mathbf{x}}{\|\mathbf{x}\|}, \frac{\mathbf{z}}{\|\mathbf{z}\|} \right\rangle$
 - popular choice for measuring similarity of text documents
 - normalizing (dividing by L_2 -norm) projects vectors onto unit sphere, their dot product is the cosine of the angle between the vectors
- many more ...

Based on slide by Eric Eaton

Kernel Construction

Instead of determining whether $k(\mathbf{x}, \mathbf{z})$ is a valid kernel, construct $k(\mathbf{x}, \mathbf{z})$ from simpler kernels (active area of ML).

Closure Properties of Kernels

Let $k_1(\mathbf{x}, \mathbf{z})$ and $k_2(\mathbf{x}, \mathbf{z})$ be valid kernels with feature mappings $\phi^{(1)}(\mathbf{x})$ and $\phi^{(2)}(\mathbf{x})$. Then

- (addition) $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$
 - (scaling) $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) k_1(\mathbf{x}, \mathbf{z}) f(\mathbf{z})$
for any real-valued function $f: \mathbb{R}^d \mapsto \mathbb{R}$
 - (multiplication) $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) k_2(\mathbf{x}, \mathbf{z})$
- are all valid kernels.

Based on notes by Tommi Jaakola

Proofs

Two options:

- (1) Determine implicit feature mapping $\phi(\mathbf{x})$.
- (2) Show that kernel matrix is symmetric PSD.

Addition: $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z}) + k_2(\mathbf{x}, \mathbf{z})$

Based on notes by Tommi Jaakola



Prove that RBF Kernel is a Valid Kernel

Assume $\sigma^2 = 1$.

(Easy to generalize, or recognize $ak(\mathbf{x}, \mathbf{z})$ for $a > 0$ is a valid kernel.)

Based on notes by Tommi Jaakola



SVMs

Learning Goals

- ✓ Solving the SVM optimization problem
 - using Lagrange multipliers (leads to dual problem)
- ✓ Allowing misclassified examples
 - using slack variables (leads to soft-margin SVM)
- ✓ Describe the SVM loss function
- ✓ Allowing non-linear decision boundaries
 - using kernels

Take-Aways

- Maximum-margin separator
- Primal-dual formulation
- Hard vs soft-margin SVM
- Hinge loss
- Kernels (“kernel trick”)

SVM Practical Advice

When Applying SVMs

Use SVM software package to solve for parameters

- e.g. SVMlight, libsvm, cvx (fast!), etc

Need to specify

- Choice of hyperparameter C
- Choice of kernel function
- Associated kernel parameters
e.g. p for polynomial kernel, σ for RBF kernel

Based on slide by Eric Eaton

Practical Advice

When faced with ML problem, it is sometimes not clear which algorithm to use.

The algorithm matters, but what matters more are

- How much data you have
- How good you are at error analysis and debugging learning algorithms
- How you design features
- (Upcoming lecture: Advice for Applying ML)