

## Introduction to Kernel Methods

Ron Parr  
CPS 271

## Outline

- Motivation for kernel methods
- The dual view
- What makes a kernel?
- The Gaussian Process view
- Review of why we care
- Classification

## Parametric Methods

- Methods discussed so far are *parametric*
- Learning mechanism/representation is characterized by a fixed set of parameters
  - Regression coefficients
  - Neural network weights
  - Linear discriminant means and covariances
- Connection between data and output is complicated and circuitous:
  - First we pick an algorithm, then our features,
  - Algorithm tunes parameters to data
  - Data are discarded
- Final result is a function from data to predictions
- Is there a more direct route?

## Non-Parametric Methods

- Traditionally the simplest methods
- Not discussed heavily in the book
- K-Nearest neighbor
  - Classify according to nearest (in feature space) neighbor
  - For  $k > 1$ , vote
- Locally weighted regression
  - Use query dependent regression weights (rerun regression for every query)
  - Weights decay with distance from query point
- Issues
  - Defining closeness (both scale and dimensions)
  - Efficiency
  - Tendency to overfit w/o careful tweaking

## But how different are these, really?

- Both are functions from data to predictions
- Non-parametric methods
  - Seem more constrained
  - Store less information
- Is there a unifying view?

## Dual View of Regression

- Recall (regularized) regression minimizes:

$$J(\mathbf{w}) = 0.5\lambda\mathbf{w}^T\mathbf{w} + 0.5\sum_{i=1}^N (\mathbf{w}^T\phi(x^{(i)}) - t^{(i)})^2$$

- Setting  $dJ/d\mathbf{w} = 0$ :

$$\mathbf{w} = \frac{-1}{\lambda} \sum_{i=1}^N (\mathbf{w}^T\phi(x^{(i)}) - t^{(i)})\phi(x^{(i)})$$

### Dual View Continued

$$\mathbf{w} = \frac{-1}{\lambda} \sum_{i=1}^N (\mathbf{w}^T \phi(x^{(i)}) - t^{(i)}) \phi(x^{(i)})$$

- Some notation:
 
$$\Phi = \begin{bmatrix} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \vdots \\ \phi(x^{(N)}) \end{bmatrix} \quad \mathbf{a} = \begin{bmatrix} \frac{-1}{\lambda} (\mathbf{w}^T \phi(x^{(1)}) - t^{(1)}) \\ \frac{-1}{\lambda} (\mathbf{w}^T \phi(x^{(2)}) - t^{(2)}) \\ \vdots \\ \frac{-1}{\lambda} (\mathbf{w}^T \phi(x^{(N)}) - t^{(N)}) \end{bmatrix}$$
- Simplifying:  $\mathbf{w} = \Phi^T \mathbf{a}$

### Substituting

$$J(\mathbf{w}) = 0.5 \lambda \mathbf{w}^T \mathbf{w} + 0.5 \sum_{i=1}^N (\mathbf{w}^T \phi(x^{(i)}) - t^{(i)})^2$$

$$\mathbf{w} = \Phi^T \mathbf{a}$$

$$J(\mathbf{a}) = 0.5 \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + 0.5 \mathbf{t}^T \mathbf{t} - \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

Comment: This is a mess, but it makes the dependence on  $\Phi \Phi^T$  clearer

### The Gram Matrix

- Define  $K = \Phi \Phi^T$

$$J(\mathbf{a}) = 0.5 \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + 0.5 \mathbf{t}^T \mathbf{t} - \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

$$= 0.5 \mathbf{a}^T K K \mathbf{a} - \mathbf{a}^T K \mathbf{t} + 0.5 \mathbf{t}^T \mathbf{t} - \frac{\lambda}{2} \mathbf{a}^T K \mathbf{a}$$

- Setting  $dJ/d\mathbf{a} = 0$

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

### Prediction

- Recall that to predict a new value we use:
 
$$y(x) = \mathbf{w}^T \phi(x)$$
- Substituting:
 
$$\mathbf{w} = \Phi^T \mathbf{a} \quad \mathbf{a} = (K + \lambda I)^{-1} \mathbf{t}$$

$$y(x) = \mathbf{w}^T \phi(x) = \mathbf{a}^T \Phi \phi(x) = \underbrace{k(x)}_{\text{Extension to the Gram matrix}}^T (K + \lambda I)^{-1} \mathbf{t}$$

### Why this matters?

- We have expressed regression entirely in terms of the gram matrix,  $K$ , which we call the kernel
- But  $K$  is defined entirely of dot products between vectors in our training set
- Let's think about what  $K$  really means...

### What does dot product mean?

- Elements of  $K$  are dot products between vectors,  $x^T y$
- If  $x=y$ ,  $x^T y$  is squared magnitude
- If  $x$  has norm 1, then  $x^T y$  is projection of  $y$  onto  $x$
- In general  $x^T y / (||x|| ||y||)$  is the cosine between  $x$  and  $y$
- If  $x$  and  $y$  have the same magnitude, then  $x^T y$  is maximized when  $x=y$

### What happens if we redefine dot product?

- If we change the meaning of  $x^T y$ , then we change what it means to two vectors to be similar
- The big idea: Turn feature on its head
- Instead of invested effort in finding features, invest effort in finding kernels
- Algorithmically equivalent to redefining the K matrix

### Kernel example

- The default (linear) kernel says  $k(x,y) = x^T y$
- What if  $k(x,y) = (x^T y)^2$  instead?
- Suppose  $x = (x_1, x_2)$ ,  $y = (y_1, y_2)$

$$\begin{aligned}
 k(x,y) &= x^T y = (x_1 y_1 + x_2 y_2)^2 \\
 &= (x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2) \\
 &= \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}^T \begin{pmatrix} y_1^2 & \sqrt{2} y_1 y_2 & y_2^2 \end{pmatrix} \\
 &= \phi(x)^T \phi(y)
 \end{aligned}
 \quad \phi(z) = \begin{pmatrix} z_1^2 \\ \sqrt{2} z_1 z_2 \\ z_2^2 \end{pmatrix}^T$$

### What just happened?

- We just changed the feature space from linear in  $x$ , to quadratic in  $x$
- But we don't actually need to construct these features; we just need to redefine  $K$
- Question: Does every choice of  $K$  make sense? (Does every  $K$  correspond to sensible basis?)

### What makes a valid kernel?

- In general, we want our  $K$  matrix to be symmetric, positive semidefinite (not proved here)
- A *sufficient* (but not necessary) condition is for  $K$  to behave like a distance metric
  - Nonnegative
  - $K(x,x) = 0$
  - Symmetric
  - Obeys triangle inequality
- Fancy kernels can be constructed by combining simple ones

### Valid Kernel Combinations

$$\begin{aligned}
 k(x, x') &= c k_1(x, x') \\
 k(x, x') &= f(x) k_1(x, x') f(x') \\
 k(x, x') &= q(k_1(x, x')) \\
 k(x, x') &= e^{k_1(x, x')} \\
 k(x, x') &= k_1(x, x') + k_2(x, x') \\
 k(x, x') &= k_1(x, x') k_2(x, x') \\
 k(x, x') &= k_3(\phi(x), \phi(x')) \\
 k(x, x') &= x^T A x' \\
 k(x, x') &= k_a(x_a, x_a') + k_b(x_b, x_b') \\
 k(x, x') &= k_a(x_a, x_a') k_b(x_b, x_b')
 \end{aligned}$$

Assumptions:  
 $c > 0$   
 $f(\cdot)$  is any function  
 $q(\cdot)$  is polynomial w/  $\geq 0$  coefficients  
 $k_i$  is any valid kernel function  
 $x = (x_a, x_b)$   
 $A$  is symmetric and PSD

### Other Interesting Kernels

- Polynomial kernel:
 
$$k(x, x') = (x^T x' + c)^M$$
- Gaussian Kernel:
 
$$k(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$$

## More Kernels

- Kernels over sets:

$$k(A, A') = 2^{|A \cap A'|}$$

- Kernels from probabilities

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

## Gaussian Process View

- Recall that regularized regression is equivalent to the ML solution to Bayesian regression with a mean 0 Gaussian prior on the weights
- In the dual view, we interpret our training data as inducing a joint distribution over the y targets

## Gaussian Process Formulation

$$\mathbf{y} = \Phi \mathbf{w}$$

- $\mathbf{w}$  is a vector of Gaussian rvs, so  $\mathbf{y}$  is a vector of Gaussian rvs

$$E(\mathbf{y}) = \Phi E(\mathbf{w}) = \mathbf{0}$$

$$\text{cov}(\mathbf{y}) = E(\mathbf{y}\mathbf{y}^T) = \Phi E(\mathbf{w}\mathbf{w}^T)\Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K}$$

- Notes:
  - Prior variance on  $\mathbf{w} = \alpha$
  - We hid this in  $\mathbf{K}$

## Noise model for targets

- Recall that we assume our targets values  $\mathbf{t}$ , are the result of an underlying  $\mathbf{y}$ , corrupted with Gaussian noise

$$p(\mathbf{y}) = N(\mathbf{y} | \mathbf{0}, \mathbf{K})$$

$$p(\mathbf{t} | \mathbf{y}) = N(\mathbf{t} | \mathbf{y}, \beta^{-1})$$

$$p(\mathbf{t} | \mathbf{y}) = N(\mathbf{t} | \mathbf{y}, \beta^{-1} \mathbf{I})$$

$$p(\mathbf{t}) = \int p(\mathbf{t} | \mathbf{y}) p(\mathbf{y}) d\mathbf{y} = N(\mathbf{t} | \mathbf{0}, \mathbf{C})$$

$$\mathbf{C} = \mathbf{K} + \beta^{-1} \mathbf{I}$$

## Extending to Prediction

- For prediction, we want

$$p(\mathbf{t}_{N+1}) = N(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

- We don't need  $\mathbf{y}$  to define  $\mathbf{C}_{N+1}$

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}$$

- Where  $\mathbf{C}_N$  is our previous  $\mathbf{C}$ ,  $\mathbf{k}$  is the kernel between old (training) and new ( $\mathbf{x}_{N+1}$ ) points, and  $c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$

## Now we turn the magic crank

- All we need to do now is to condition on previously observed target values to get

$$P(\mathbf{t}_{N+1} | \mathbf{t})$$

- Using the conditioning results from chapter 2:

$$\mu(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$$

## Compare with Regularized Regression

- Regularized regression with kernels:

$$y(x) = \mathbf{w}^T \phi(x) = \mathbf{a}^T \Phi \phi(x) = \mathbf{k}(x)^T (K + \lambda I)^{-1} \mathbf{t}$$

- Gaussian process regression:

$$\begin{aligned} \mu(x_{N+1}) &= \mathbf{k}^T C_N^{-1} \mathbf{t} \\ \sigma^2(x_{N+1}) &= c - \mathbf{k}^T C_N^{-1} \mathbf{k} \end{aligned}$$

- Recall:

$$C = \mathbf{K} + \beta^{-1} \mathbf{I}$$

## Advantages of the GP View

- Why bother with GP if regularized regression does the same thing?
- GP also gives us a variance in prediction
- GP gives us a distribution over targets
- GP is more general, can incorporate other types of priors, including priors over priors

## Reality Check

- Cost of ordinary regression:
  - Feature space of dimension  $k$ ,  $N$  training points
  - Storage of solution  $O(k)$
  - Computation:
    - Cubic in  $k$
    - linear in  $N$
- Cost of kernel version:
  - Feature space of dimension  $k$ ,  $N$  training points
  - Storage of solution  $O(N)$
  - Computation
    - Cubic in  $N$ ,
    - No explicit dependence on  $k$

## Why do we like Kernels?

- Let us experiment with feature spaces without paying the cost of constructing the features
- But what about overfitting? (Isn't  $k > N$  dangerous?)
- Yes! This is why we need regularization!
- This issue becomes particularly interesting in the context of support vector machines

## Kernels for Classification

- Idea 1:
  - Use logistic regression, replacing  $\mathbf{x}^T \mathbf{w}$  with some kind of kernel regression
  - Problem:
    - No clean training algorithm
    - Must use approximations
- Idea 2: Support Vector Machines (next chapter)

## Relationship to Neural Networks

- Assume:
  - Linear output node
  - Large number of hidden nodes
  - Gaussian prior on weights
- Can show that in the limit of a large number of hidden nodes, the above neural network behaves like a Gaussian process

## Some Concluding Thoughts

- Kernels and Gaussian processes have become *very* popular in recent years
- Why?
  - Coolness factor
  - Ability to work in weird spaces implicitly
- Is this a good thing?
  - Not a substitute for standard linear + good features
  - Some advantages over other non-parametric methods
    - Potential for elegant treatment of regularization
    - GP provides probability distribution
  - Kernels useful for classification and other techniques beyond regression