

# MACHINE LEARNING AND PATTERN RECOGNITION:

## Lecture 3.1:

### Basis Functions, Kernel Methods

Yann LeCun,  
Courant Institute, NYU

# Linear Machines: Regression with Mean Square

---

## Linear Regression, Mean Square Loss:

- decision rule:  $y = W'X$
- loss function:  $L(W, y^i, X^i) = \frac{1}{2}(y^i - W'X^i)^2$
- gradient of loss:  $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(y^i - W(t)'X^i)X^i$
- update rule:  $W(t+1) = W(t) + \eta(t)(y^i - W(t)'X^i)X^i$
- direct solution: solve linear system  $[\sum_{i=1}^P X^i X^{i'}]W = \sum_{i=1}^P y^i X^i$

# Linear Machines: Perceptron

---

## Perceptron:

- decision rule:  $y = F(W'X)$  ( $F$  is the threshold function)
- loss function:  $L(W, y^i, X^i) = (F(W'X^i) - y^i)W'X^i$
- gradient of loss:  $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(y^i - F(W'X^i))X^i$
- update rule:  $W(t+1) = W(t) + \eta(t)(y^i - F(W'X^i))X^i$
- direct solution: find  $W$  such that  $-y^i F(W'X^i) < 0 \quad \forall i$

# Linear Machines: Logistic Regression

---

## Logistic Regression, Negative Log-Likelihood Loss function:

- decision rule:  $y = F(W'X)$ , with  $F(a) = \frac{1 - \exp(a)}{1 + \exp(a)}$  (sigmoid function).
- loss function:  $L(W, y^i, X^i) = 2 \log(1 + \exp(-y^i W'X^i))$
- gradient of loss:  $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(Y^i - F(W'X)) X^i$
- update rule:  $W(t + 1) = W(t) + \eta(t)(y^i - F(W(t)'X^i))X^i$

# General Gradient-Based Supervised Learning Machine

---

## Neural Nets, and many other models:

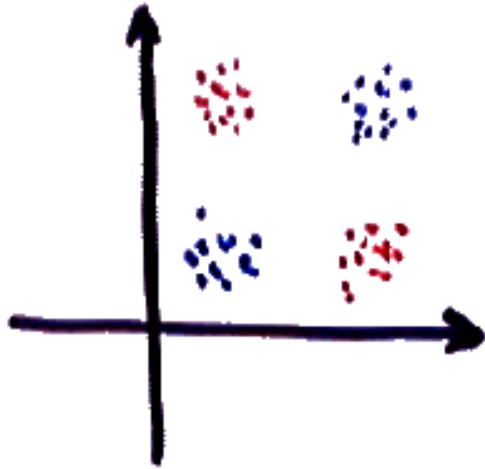
- decision rule:  $y = F(W, X)$ , where  $F$  is some function, and  $W$  some parameter vector.
- loss function:  $L(W, y^i, X^i) = D(y^i, F(W, X^i))$ , where  $D(y, f)$  measures the “discrepancy” between  $A$  and  $B$ .
- gradient of loss:  $\frac{\partial L(W, y^i, X^i)}{\partial W} = \frac{\partial D(y^i, f)}{\partial f} \frac{\partial F(W, X^i)}{\partial W}$
- update rule:  $W(t + 1) = W(t) - \eta(t) \frac{\partial D(y^i, f)}{\partial f} \frac{\partial F(W, X^i)}{\partial W}$

## Three Questions:

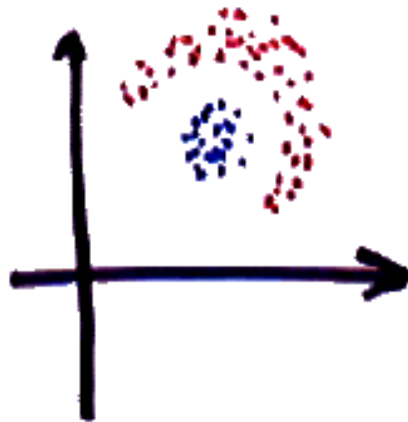
- What architecture  $F(W, X)$ .
- What loss Function  $L(W, y^i, X^i)$ .
- What optimization method.

# Limitations of Linear Machines

---

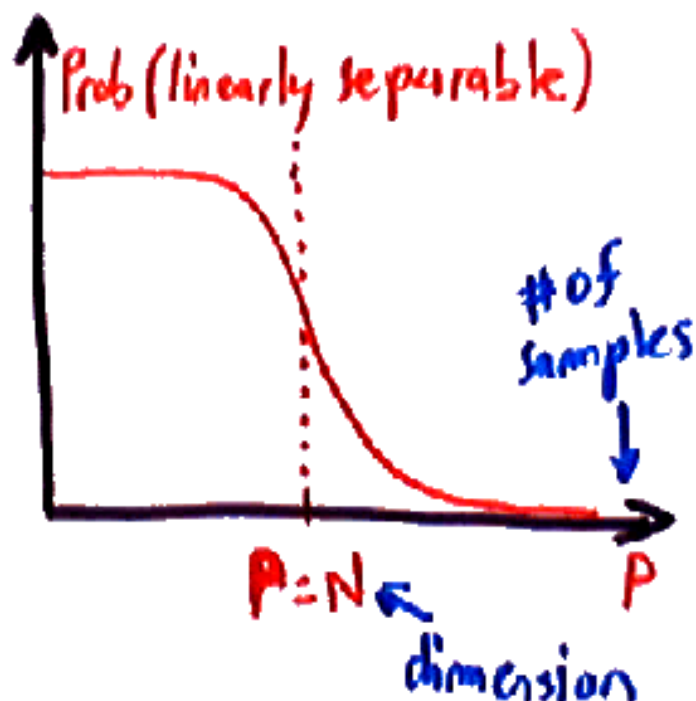


The *Linearly separable* dichotomies are the partitions that are realizable by a linear classifier (the boundary between the classes is a hyperplane).



# Number of Linearly Separable Dichotomies

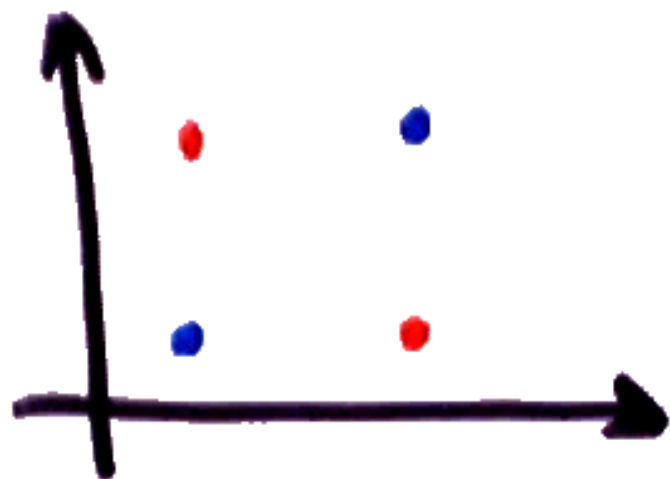
The probability that  $P$  samples of dimension  $N$  are linearly separable goes to zero very quickly as  $P$  grows larger than  $N$  (Cover's theorem, 1966).



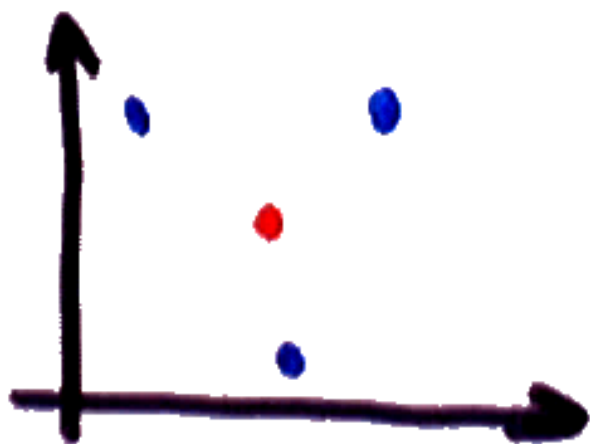
- Problem: there are  $2^P$  possible dichotomies of  $P$  points.
- Only about  $N$  are linearly separable.
- If  $P$  is larger than  $N$ , the probability that a random dichotomy is linearly separable is very, very small.

# Example of Non-Linearly Separable Dichotomies

---

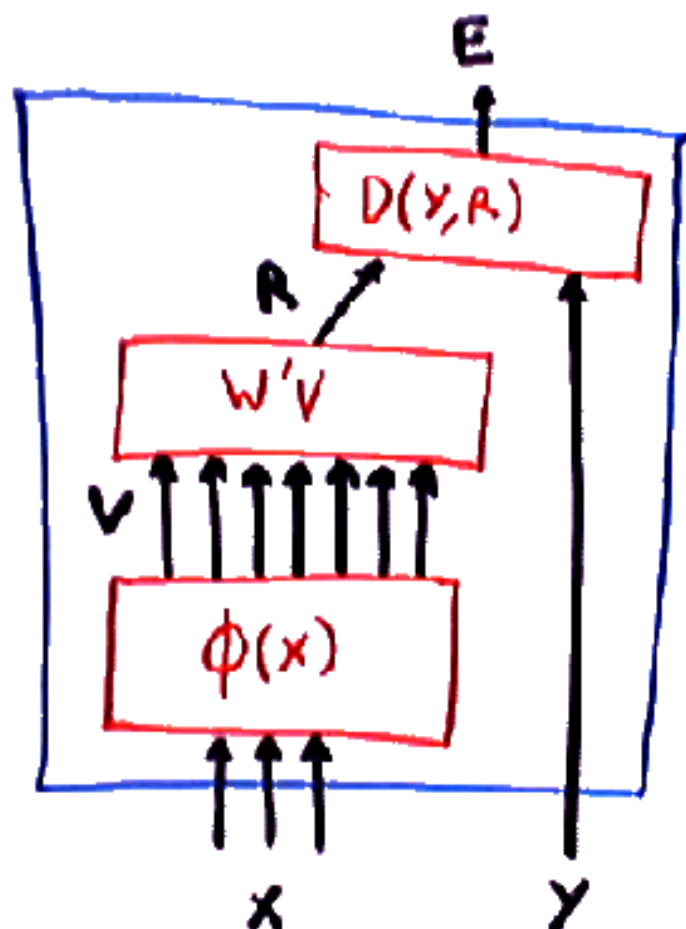


- Some seemingly simple dichotomies are not linearly separable
- **Question:** How do we make a given problem linearly separable?



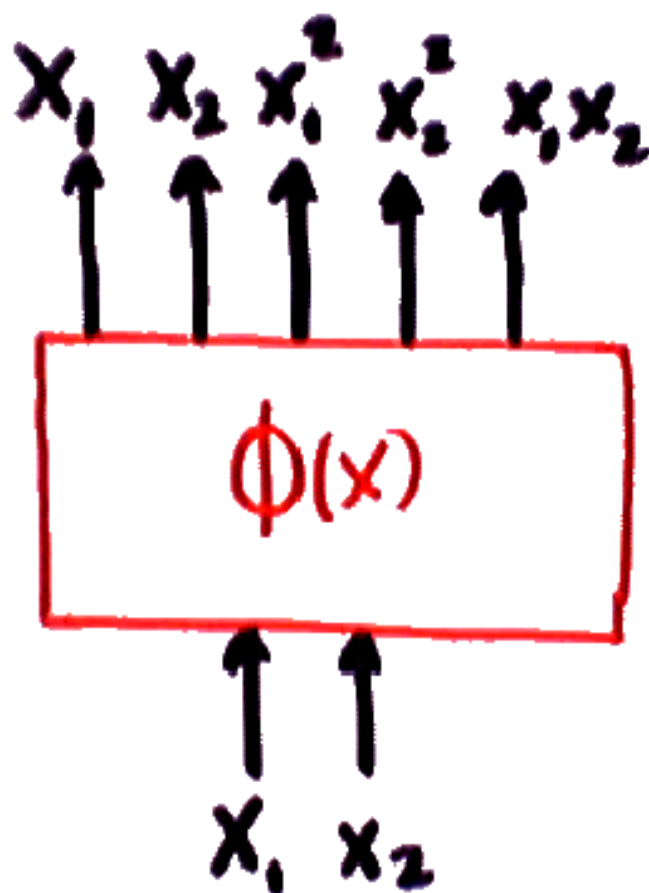


# Making $N$ Larger: Preprocessing



- **Answer 1:** we make  $N$  larger by augmenting the input variables with new “features”.
- we map/project  $X$  from its original  $N$ -dimensional space into a higher dimensional space where things are more likely to be linearly separable, using a vector function  $\Phi(X)$ .
- $E(Y, X, W) = D(Y, R)$
- $R = f(W'V)$
- $V = \Phi(X)$

# Adding Cross-Product Terms



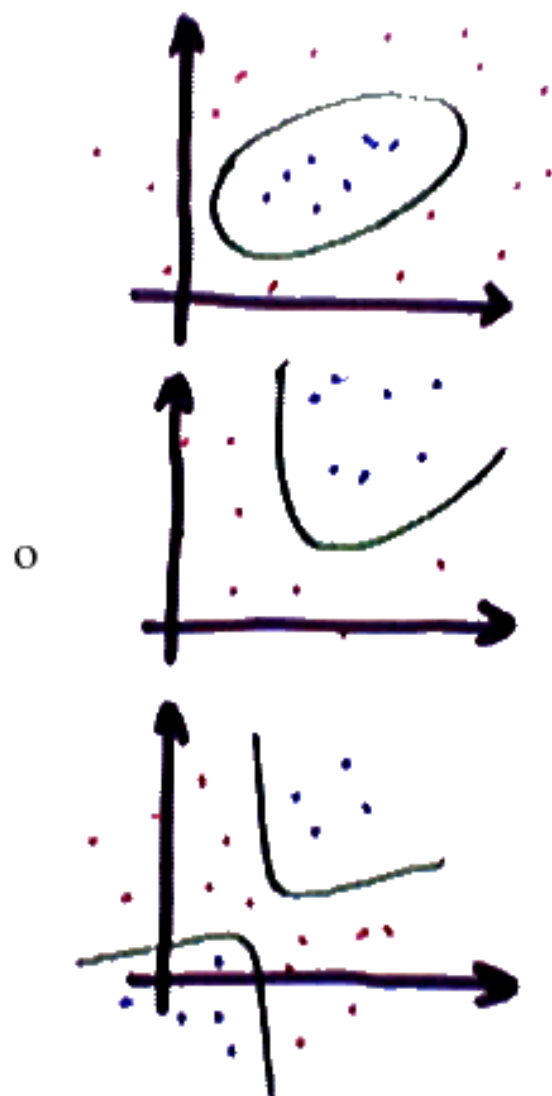
- Polynomial Expansion.
- If our original input variables are  $(1, x_1, x_2)$ , we construct a new *feature vector* with the following components:

$$\Phi(1, x_1, x_2) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$$

i.e. we add all the cross-products of the original variables.

- we map/project  $X$  from its original  $N$ -dimensional space into a higher dimensional space with  $N(N+1)/2$  dimensions.

# Polynomial Mapping



- Many new functions are now separable with the new architecture.
- With cross-product features, the family of class boundaries in the original space is the conic sections (ellipse, parabola, hyperbola).
- to each possible boundary in the original space corresponds a linear boundary in the transformed space.
- Because this is essentially a linear classifier with a preprocessing, we can use standard linear learning algorithms (perceptron, linear regression, logistic regression...).

# Problems with Polynomial Mapping

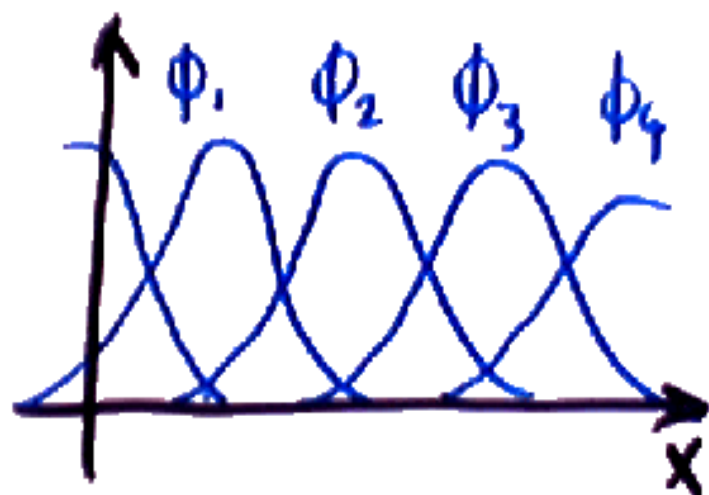
---

- We can generalize this idea to higher degree polynomials, adding cross-product terms with 3, 4 or more variables.
- Unfortunately, the number of terms is the number of combinations  $d$  choose  $N$ , which grows like  $N^d$ , where  $d$  is the degree, and  $N$  the number of original variables.
- In particular, the number of free parameters that must be learned is also of order  $N^d$ .
- This is impractical for large  $N$  and for  $d > 2$ .
- Example: handwritten digit recognition (16x16 pixel images). Number of variables: 256. Degree 2: 32,896 variables. Degree 3: 2,796,160. Degree 4: 247,460,160.....

## Next Idea: Tile the Space

---

place a number of equally-spaced “bumps” that cover the entire input space.



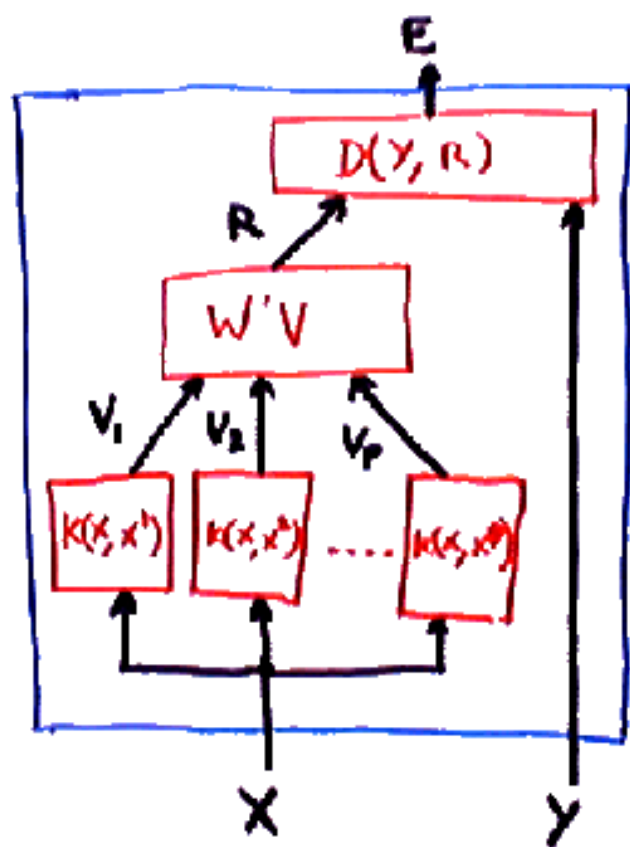
- For classification, the bumps can be Gaussians
- For regression, the basis functions can be wavelets, sine/cosine, splines (pieces of polynomials)....
- **problem:** this does not work with more than a few dimensions.
- The number of bumps necessary to cover an  $N$  dimensional space grows exponentially with  $N$ .

# Sample-Centered Basis Functions (Kernels)

Place the center of a basis function around each training sample. That way, we only spend resources on regions of the space where we actually have training samples.

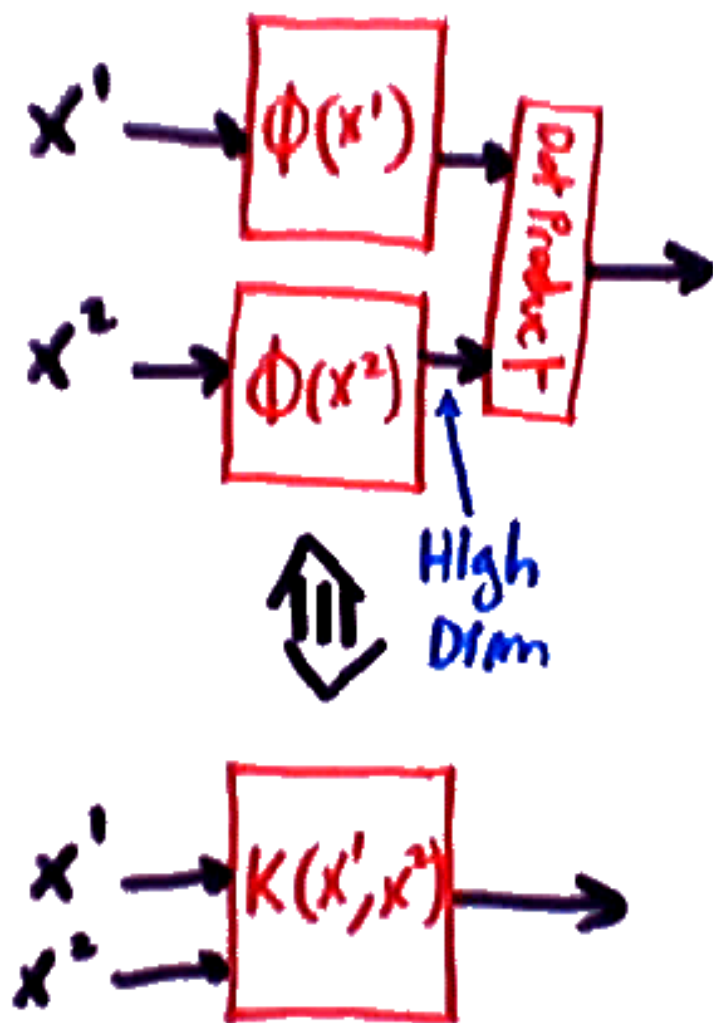
- Discriminant function:

$$f(X, W) = \sum_{k=1}^{k=P} W_k K(X, X^k)$$



- $K(X, X')$  often takes the form of a *radial basis function*:  
 $K(X, X') = \exp(b\|X - X'\|^2)$  or a polynomial  $K(X, X') = (X \cdot X' + 1)^m$
- This is a very common architecture, which can be used with a number of energy functions.
- In particular, this is the architecture of the so-called **Support Vector Machine** (SVM), but the energy function of the SVM is a bit special. We will study it later in the course.

# The Kernel Trick



- If the kernel function  $K(X, X')$  verifies the *Mercer conditions*, then there exist a mapping  $\Phi$ , such that  $\Phi(X) \cdot \Phi(X') = K(X, X')$ .
- The Mercer conditions are that  $K$  must be symmetric, and must be positive definite (i.e  $K(X, X)$  must be positive for all  $X$ ).
- In other words, if we want to map our  $X$  into a high-dimensional space (so as to make them linearly separable), and all we have to do in that space is compute dot products, we can take a shortcut and simply compute  $K(X^1, X^2)$  without going through the high-dimensional space.
- This is called the “**kernel trick**”. It is used in many so-called Kernel-based methods, including Support Vector Machines.

# Examples of Kernels

---

- **Quadratic kernel:**  $\Phi(X) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2)$  then

$$K(X, X') = \Phi(X) \cdot \Phi(X') = (X \cdot X' + 1)^2$$

- **Polynomial kernel:** this generalizes to any degree  $d$ . The kernel that corresponds to  $\Phi(X)$  being a polynomial of degree  $d$  is

$$K(X, X') = \Phi(X) \cdot \Phi(X') = (X \cdot X' + 1)^d.$$

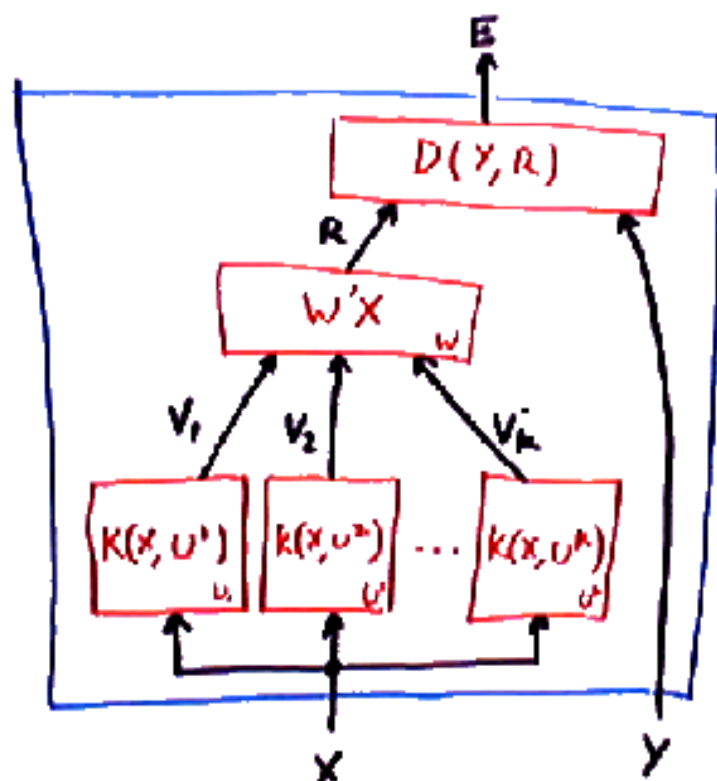
- **Gaussian Kernel:**

$$K(X, X') = \exp(-b\|X - X'\|^2)$$

This kernel, sometimes called the Gaussian Radial Basis Function, is very commonly used.



# Sparse Basis Functions

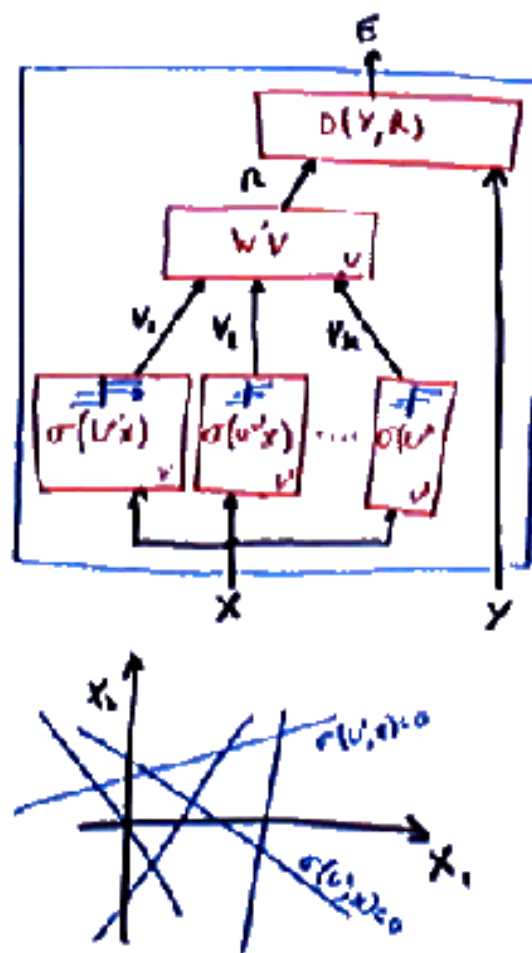


- Place the center of a basis function around areas containing training samples.
- Idea 1: use an unsupervised clustering algorithm (such as K-means or mixture of Gaussians) to place the centers of the basis functions in areas of high sample density.
- Idea 2: adjust the basis function centers through gradient descent in the loss function.

The discriminant function  $F$  is:

$$F(X, W, U^1, \dots, U^K) = \sum_{k=1}^{k=K} W_k K(X, U^k)$$

## Other Idea: Random Directions



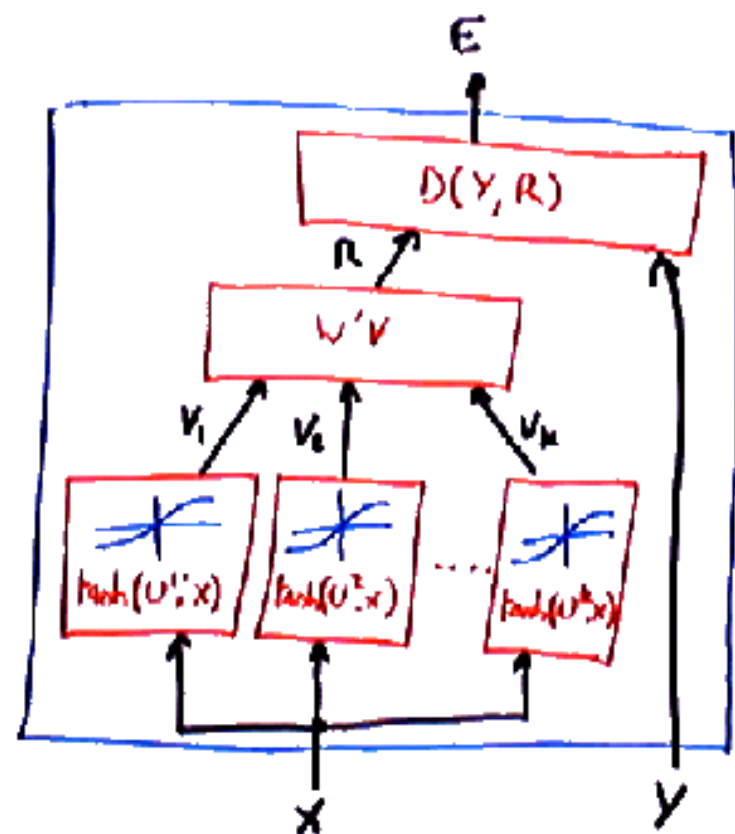
- Partition the space in lots of little domains by randomly placing lots of hyperplanes.
- Use many variables of the type  $q(W^k X)$ , where  $q$  is the threshold function (or some other squashing function) and  $W_k$  is a randomly picked vector.
- This is the original Perceptron.
- Without the non-linearity, the whole system would be linear (product of linear operations), and therefore would be no more powerful than a linear classifier.
- **problem:** a bit of a wishful thinking, but it works occasionally.

# Neural Net with a Single Hidden Layer

A particularly interesting type of basis function is the sigmoid unit:  $V_k = \tanh(U'^k X)$

- a network using these basis functions, whose output is  $R = \sum_{k=1}^{k=K} W_k V_k$  is called a *single hidden-layer neural network*.
- Similarly to the RBF network, we can compute the gradient of the loss function with respect to the  $U^k$ :

$$\begin{aligned}\frac{\partial L(W)}{\partial U^j} &= \frac{\partial L(W)}{\partial R} W_j \frac{\partial \tanh(U'_j X)}{\partial U_j} \\ &= \frac{\partial L(W)}{\partial R} W_j \tanh'(U'_j X) X'\end{aligned}$$



Any well-behaved function can be approximated as close as we wish by such networks (but  $K$  might be very large).