# Support Vector Machines

# Linear Classifiers



$$\alpha$$

$$\mathbf{x} \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$$f(\mathbf{x}, \mathbf{w}, b) = sign(\mathbf{w} \; \mathbf{x} + b)$$

- ● denotes +1
- ○ denotes -1

$$\mathbf{w} \; \mathbf{x} + b > 0$$

$$\mathbf{w} \; \mathbf{x} + b = 0$$

How would you classify this data?

$$\mathbf{w} \; \mathbf{x} + b < 0$$

# How do we characterize "power"?

- Different machines have different amounts of "power".

- Tradeoff between:
  - More power: Can model more complex classifiers but might overfit.
  - Less power: Not going to overfit, but restricted in what it can model.

- How do we characterize the amount of power?

# Some definitions

- Given some machine **f**
- And under the assumption that all training points $(x_k, y_k)$ were drawn i.i.d from some distribution.
- And under the assumption that future test points will be drawn from the same distribution
- Define

$$R(\alpha) = \text{TESTERR}(\alpha) = E\left[\frac{1}{2}|y - f(x,\alpha)|\right] = \begin{array}{l}\text{Probability of} \\ \text{Misclassification}\end{array}$$

Official terminology

Terminology we'll use

# Some definitions

- Given some machine **f**
- And under the assumption that all training points $(x_k, y_k)$ were drawn i.i.d from some distribution.
- And under the assumption that future test points will be drawn from the same distribution
- Define

$$R(\alpha) = \text{TESTERR}(\alpha) = E\left[\frac{1}{2}|y - f(x, \alpha)|\right] = \begin{array}{l} \text{Probability of} \\ \text{Misclassification} \end{array}$$

Official terminology

Terminology we'll use

$$R^{emp}(\alpha) = \text{TRAINERR}(\alpha) = \frac{1}{R}\sum_{k=1}^{R}\frac{1}{2}|y_k - f(x_k, \alpha)| = \begin{array}{l} \text{Fraction Training} \\ \text{Set misclassified} \end{array}$$

R = #training set data points

# Vapnik-Chervonenkis dimension

$$\text{TESTERR}(\alpha) = E\left[\frac{1}{2}\left|y - f(x,\alpha)\right|\right] \qquad \text{TRAINERR}(\alpha) = \frac{1}{R}\sum_{k=1}^{R}\frac{1}{2}\left|y_k - f(x_k,\alpha)\right|$$

- Given some machine **f**, let *h* be its VC dimension.
- *h* is a measure of **f**'s power (*h* does not depend on the choice of training set)
- Vapnik showed that with probability 1-η

$$\text{TESTERR}(\alpha) \leq \text{TRAINERR}(\alpha) + \sqrt{\frac{h(\log(2R/h)+1) - \log(\eta/4)}{R}}$$

This gives us a way to estimate the error on future data based only on the training error and the VC-dimension of *f*

# What VC-dimension is used for

$$\text{TESTERR}(\alpha) = E\left[\frac{1}{2}|y - f(x, \alpha)|\right] \qquad \text{TRAINERR}(\alpha) = \frac{1}{R}\sum_{k=1}^{R}\frac{1}{2}|y_k - f(x_k, \alpha)|$$

- Given some machine **f**, let $h$ be its VC dimension.
- $h$ is a measure of **f**'s power.
- Vapnik showed that with probability $1 - \eta$:

$$\text{TESTERR}(\alpha) \le \text{TRAINERR}(\alpha) + \sqrt{\frac{h(\log(2R/h) + 1) - \log(\eta/4)}{R}}$$

This gives us a way to estimate the error on future data based only on the training error and the VC-dimension of f

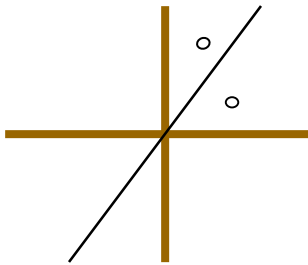**But given machine f, how do we define and compute h?**

# Shattering

- Machine f can *shatter* a set of points $x_1, x_2 .. x_r$ if and only if…

  For every possible training set of the form $(x_1, y_1)$, $(x_2, y_2)$, … $(x_r, y_r)$

  …There exists some value of $\alpha$ that gets zero training error.

There are $2^r$ such training sets to consider, each with a different combination of +1's and −1's for the y's
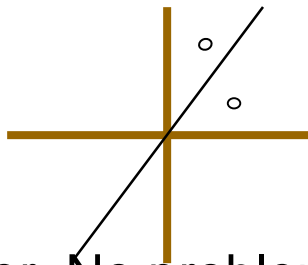
# Shattering

- Machine f can *shatter* a set of points $x_1, x_2 .. X_r$ if and only if…

    For every possible training set of the form $(x_1,y_1)$ , $(x_2,y_2)$ ,… $(x_r,y_r)$

    …There exists some value of $\alpha$ that gets zero training error.

- Question: Can the following f shatter the following points?
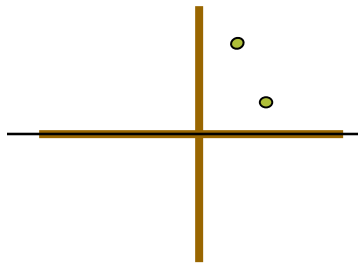
$$f(x,w) = sign(x.w)$$

# Shattering

- Machine f can *shatter* a set of points $x_1, x_2 .. X_r$ if and only if…

    For every possible training set of the form $(x_1, y_1), (x_2, y_2), … (x_r, y_r)$

    …There exists some value of $\alpha$ that gets zero training error.

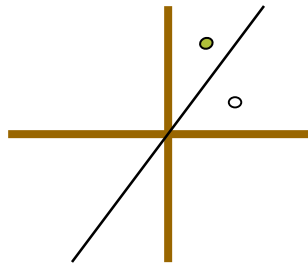- Question: Can the following f shatter the following points?

$$f(x, w) = \text{sign}(x.w)$$

- Answer: No problem. There are four training sets to consider

w=(0,1)    w=(-2,3)    w=(2,-3)    w=(0,-1)

# Shattering

- Machine f can *shatter* a set of points $x_1, x_2 .. X_r$ if and only if…

    For every possible training set of the form $(x_1, y_1), (x_2, y_2), … (x_r, y_r)$

    …There exists some value of $\alpha$ that gets zero training error.

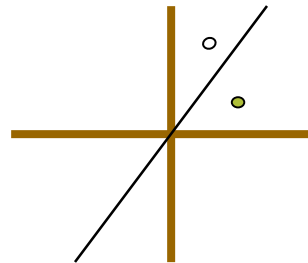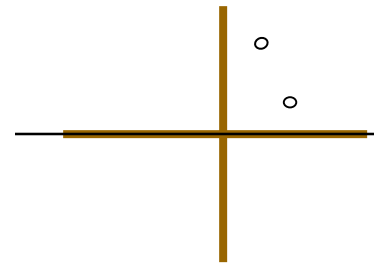- Question: Can the following f shatter the following points?

$$f(x, b) = \text{sign}(x.x - b)$$
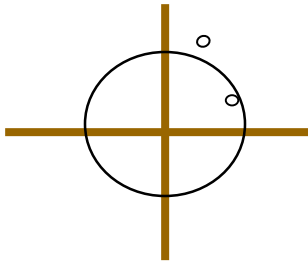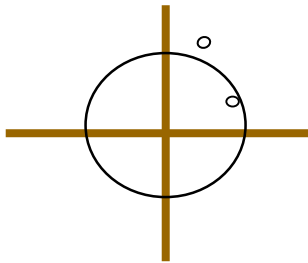
# Shattering

- Machine f can *shatter* a set of points $x_1, x_2 .. X_r$ if and only if…

   For every possible training set of the form $(x_1,y_1)$ , $(x_2,y_2)$ , … $(x_r,y_r)$

   …There exists some value of $\alpha$ that gets zero training error.

- Question: Can the following f shatter the following points?



$$f(x,b) = \text{sign}(x.x-b)$$

- Answer: No way my friend.
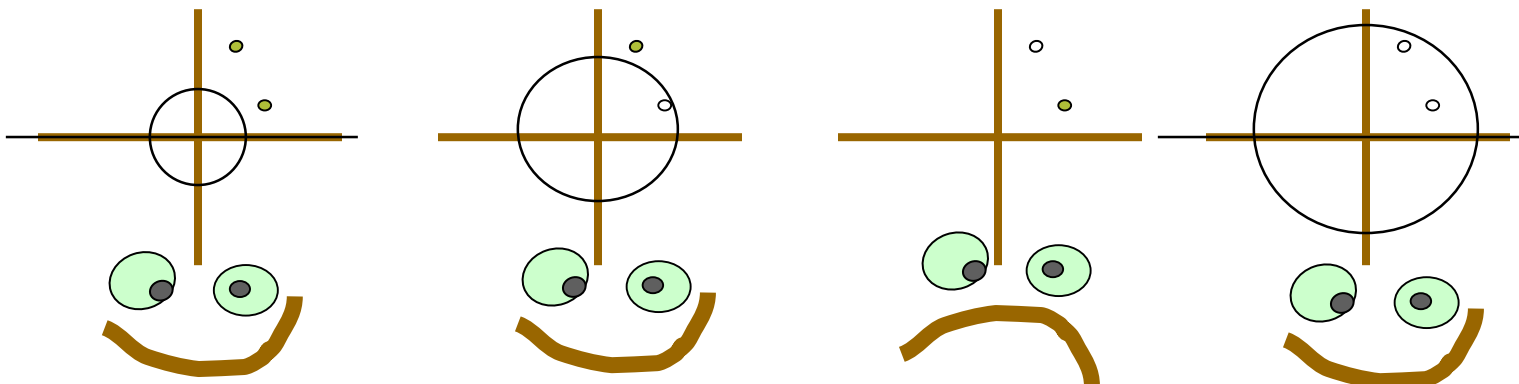
# Definition of VC dimension

Given machine **f**, the VC-dimension *h* is

The maximum number of points that can be arranged so that **f** shatter them.

Example: What's VC dimension of f(x,b) = sign(x.x-b)

# VC dim of trivial circle

Given machine *f*, the VC-dimension *h* is

The maximum number of points that can be arranged so that *f* shatter them.

Example: What's VC dimension of $f(x,b) = \text{sign}(x.x - b)$

Answer = 1: we can't even shatter two points! (but it's clear we can shatter 1)

# Reformulated circle

Given machine $f$, the VC-dimension $h$ is

> The maximum number of points that can be arranged so that $f$ shatter them.

Example: For 2-d inputs, what's VC dimension of $f(x,q,b) = \text{sign}(qx.x-b)$

# Reformulated circle

Given machine **f**, the VC-dimension *h* is

> The maximum number of points that can be arranged so that **f** shatter them.

Example: What's VC dimension of f(x,q,b) = sign(qx.x-b)

- Answer = 2

q,b are -ve

# Reformulated circle

Given machine *f*, the VC-dimension *h* is

The maximum number of points that can be arranged so that *f* shatter them.

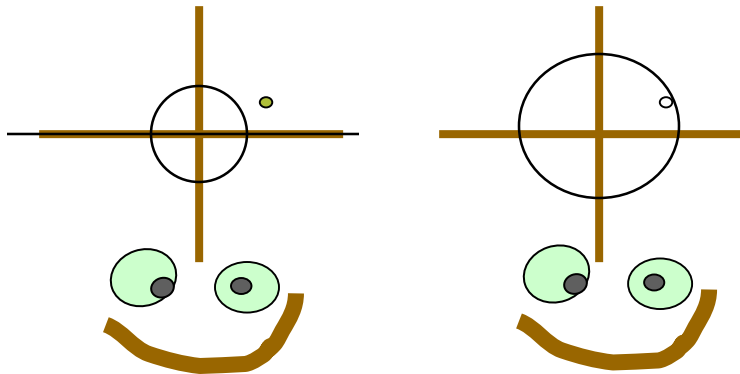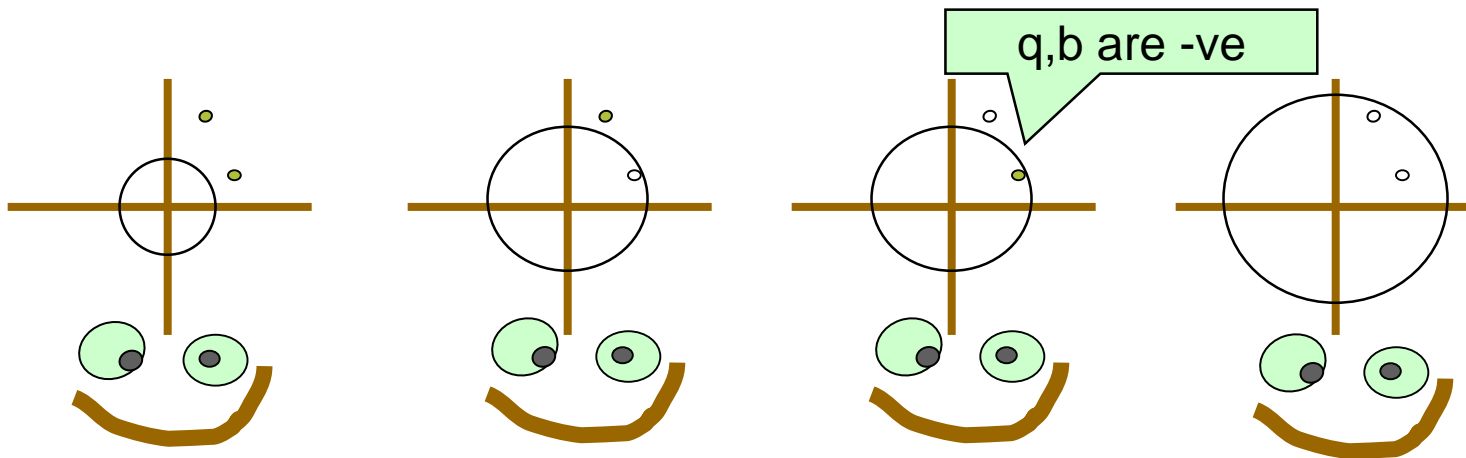Example: What's VC dimension of $f(x,q,b) = \text{sign}(qx.x-b)$

- Answer = 2 (clearly can't do 3)

q,b are -ve

# Vapnik-Chervonenkis dimension

$$\text{TESTERR}(\alpha) = E\left[\frac{1}{2}|y - f(x,\alpha)|\right] \qquad \text{TRAINERR}(\alpha) = \frac{1}{R}\sum_{k=1}^{R}\frac{1}{2}|y_k - f(x_k,\alpha)|$$

- Given some machine **f**, let $h$ be its VC dimension.
- $h$ is a measure of **f**'s power ($h$ does not depend on the choice of training set)
- Vapnik showed that with probability 1-$\eta$

$$\text{TESTERR}(\alpha) \le \text{TRAINERR}(\alpha) + \sqrt{\frac{h(\log(2R/h)+1)-\log(\eta/4)}{R}}$$

This gives us a way to estimate the error on future data based only on the training error and the VC-dimension of **f**
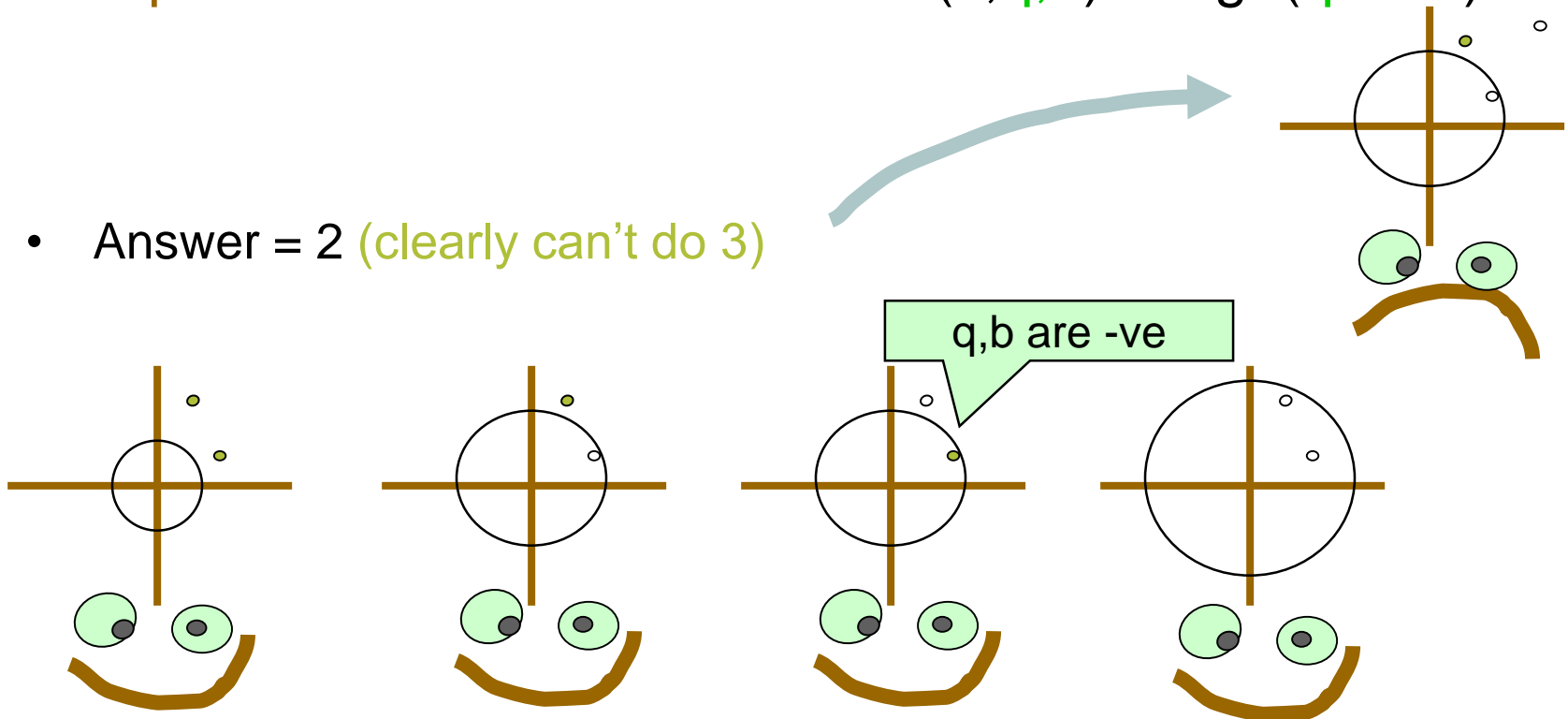
# Definition of VC dimension

Given machine **f**, the VC-dimension *h* is

The maximum number of points that can be arranged so that **f** shatter them.

Example: What's VC dimension of $f(x,b) = \text{sign}(x.x-b)$

# Linear Classifiers

$$\alpha$$

$$\boldsymbol{x} \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$$f(\boldsymbol{x}, \boldsymbol{w}, b) = sign(\boldsymbol{w}\,\boldsymbol{x} + b)$$

- denotes +1

○ denotes -1

How would you classify this data?

# Linear Classifiers

$$\alpha$$

$$\boldsymbol{x} \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$$f(\boldsymbol{x}, \boldsymbol{w}, b) = sign(\boldsymbol{w}\,\boldsymbol{x} + b)$$

- denotes +1
- denotes -1

How would you classify this data?

# Linear Classifiers

$\alpha$

$$\boldsymbol{x} \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$f(\boldsymbol{x}, \boldsymbol{w}, b) = sign(\boldsymbol{w} \ \boldsymbol{x} + b)$

- denotes +1
- denotes -1

Any of these would be fine..

..but which is best?

# Linear Classifiers

$\alpha$

$x$ → $f$ → $y^{est}$

$f(x, w, b) = sign(w\ x + b)$

- denotes +1
- denotes -1

How would you classify this data?

**Misclassified to +1 class**

# Classifier Margin

$\alpha$

$x \longrightarrow \boxed{f} \longrightarrow y^{est}$

$f(x,w,b) = sign(w\ x + b)$

- denotes +1
- denotes -1

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin

$\alpha$

$\boldsymbol{x}$ $\longrightarrow$ $f$ $\longrightarrow$ $y^{est}$

- denotes +1
- denotes -1

1. Maximizing the margin is good according to intuition
2. Implies that only support vectors are important; other training examples are ignorable.
3. Empirically it works very well.

**Support Vectors** are those datapoints that the margin pushes up against

linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

# Linear SVM Mathematically

"Predict Class = +1" zone

$x^+$

$M$=Margin Width

$x^-$

"Predict Class = -1" zone

wx+b=1
wx+b=0
wx+b=-1

What we know:
- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $w \cdot (x^+ - x^-) = 2$

$$M = \frac{(x^+ - x^-) \cdot w}{|w|} = \frac{2}{|w|}$$

# Linear SVM Mathematically

- Goal: **1) Correctly classify all training data**

$$wx_i + b \geq 1 \quad \text{if } y_i = +1$$

$$wx_i + b \leq 1 \quad \text{if } y_i = -1$$

$$y_i(wx_i + b) \geq 1 \quad \text{for all i}$$

$$M = \frac{2}{|w|}$$

**2) Maximize the Margin**

**same as minimize** $\frac{1}{2}w^t w$

- **We can formulate a Quadratic Optimization Problem and solve for w and b**

Minimize $\Phi(w) = \frac{1}{2}w^t w$

subject to $y_i(wx_i + b) \geq 1 \quad \forall i$

# Solving the Optimization Problem

Find **w** and b such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T\mathbf{w}$ is minimized;

and for all $\{(\mathbf{x_i}, y_i)\}$: $y_i(\mathbf{w}^T\mathbf{x_i} + b) \geq 1$

- **Need to optimize a *quadratic* function subject to *linear* constraints.**

- **Quadratic optimization problems are a well-known class of mathematical programming problems, and many (rather intricate) algorithms exist for solving them.**

- **The solution involves constructing a *dual problem* where a *Lagrange multiplier $\alpha_i$* is associated with every constraint in the primary problem:**

Find $\alpha_1 ... \alpha_N$ such that

$Q(\boldsymbol{\alpha}) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x_i}^T\mathbf{x_j}$ is maximized and

(1) $\Sigma\alpha_i y_i = 0$

(2) $\alpha_i \geq 0$ for all $\alpha_i$

# The Optimization Problem Solution

- The solution has the form:

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x_i} \qquad b = y_k - \mathbf{w^T x_k} \text{ for any } \mathbf{x_k} \text{ such that } \alpha_k \neq 0$$

- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x_i}$ is a support vector.

- Then the classifying function will have the form:

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x_i^T x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x_i}$.

- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x_i^T x_j}$ between all pairs of training points.

# Dataset with noise



- denotes +1
- denotes -1

- **Hard Margin: So far we require all data points be classified correctly**

  - **No training error**

- **What if the training set is noisy?**

  - **Solution 1: use very powerful kernels**

**OVERFITTING!**

# Soft Margin Classification

**Slack variables ξi can be added to allow misclassification of difficult or noisy examples.**

What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k$$

# Hard Margin v.s. Soft Margin

- **The old formulation:**

  > Find $\mathbf{w}$ and $b$ such that
  >
  > $\mathbf{\Phi}(\mathbf{w}) = \frac{1}{2}\, \mathbf{w}^{\mathrm{T}}\mathbf{w}$ is minimized and for all $\{(\mathbf{x_i}, y_i)\}$
  >
  > $y_i\,(\mathbf{w}^{\mathrm{T}}\mathbf{x_i} + \mathrm{b}) \geq 1$

- **The new formulation incorporating slack variables:**

  > Find $\mathbf{w}$ and $b$ such that
  >
  > $\mathbf{\Phi}(\mathbf{w}) = \frac{1}{2}\, \mathbf{w}^{\mathrm{T}}\mathbf{w} + C\sum \xi_i$ is minimized and for all $\{(\mathbf{x_i}, y_i)\}$
  >
  > $y_i\,(\mathbf{w}^{\mathrm{T}}\mathbf{x_i} + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$ for all $i$

- **Parameter $C$ can be viewed as a way to control overfitting.**

# Hard 1-dimensional Dataset

What would SVMs do with this data?

Not a big surprise

x=0

Positive "plane"

Negative "plane"

Doesn't look like slack variables will save us this time...

# Hard 1-dimensional Dataset

*Make up a new feature!*

Sort of…

… computed from original feature(s)

$$\mathbf{z}_k = (x_k, x_k^2)$$

Separable! MAGIC!

*x=0*

New features are sometimes called *basis functions.*

Now drop this "augmented" data into our linear SVM.

# Kernels and Linear Classifiers

Let $\vec{x} = [\vec{x}_1, \vec{x}_2] \in \mathbb{R}^2$ be a vectorial represenation
of object $x \in \mathcal{X}$

Let $\phi : \mathcal{X} \to \mathcal{K} \subset \mathbb{R}^3$ feature map be given by

$$\phi(\vec{x}) \doteq [\vec{x}_1, \vec{x}_2^2, \vec{x}_1 \vec{x}_2]^T \in \mathcal{K} \subset \mathbb{R}^3$$

**Def.** Feature space: $\mathcal{K}$

**We will use linear classifiers in this feature space.**

In the original space $\mathbb{R}^2$ for a given $\mathbf{w} \in \mathbb{R}^3$ the decision surface is:

$$\boxed{\tilde{X}_0(\mathbf{w}) = \{\vec{x} \in \mathbb{R}^2 \mid w_1 \vec{x}_1 + w_2 \vec{x}_2^2 + w_3 \vec{x}_1 \vec{x}_2 = 0\}}$$

- This is nonlinear in $\vec{x} \in \mathbb{R}^2$
- This is linear in the feature space $\phi(\vec{x}) \in \mathcal{K} \subset \mathbb{R}^3$

$$\phi(\vec{x}) \doteq [\vec{x}_1, \vec{x}_2^2, \vec{x}_1\vec{x}_2]^T \in \mathcal{K} \subset \mathbb{R}^3 \text{ feature map}$$

# Kernels and Linear Classifiers

$$\phi(\vec{x}) \doteq [\phi_1(\vec{x}), \phi_2(\vec{x}), \phi_3(\vec{x})] \doteq [\vec{x}_1, \vec{x}_2^2, \vec{x}_1\vec{x}_2]^T$$

**Feature functions**

- We seek for a small set of basis vectors $\{\phi_i\}$
  which allows perfect discrimination between
  the classes in $\mathcal{X}$ (**Feature selection**)

- If we have too many features $\Rightarrow$ overfitting can happen.

# Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?



- How about… mapping data to a higher-dimensional space:

# Non-linear SVMs: Feature spaces

■ General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: \ \mathbf{x} \rightarrow \varphi(\mathbf{x})$$

# The "Kernel Trick"

- **To produce linear separability in Higher Dimension, the linear classifier relies on dot product between vectors $K(\mathbf{x_i},\mathbf{x_j})=\mathbf{x_i}^T\mathbf{x_j}$**

- **If every data point is mapped into high-dimensional space via some transformation $\Phi:\ \mathbf{x} \rightarrow \varphi(\mathbf{x})$, the dot product becomes:**

$$K(\mathbf{x_i},\mathbf{x_j})= \varphi(\mathbf{x_i})^T\varphi(\mathbf{x_j})$$

- **A *kernel function* is some function that corresponds to an inner product in some expanded feature space.**

- **Example:**

  **2-dimensional vectors $\mathbf{x}=[x_1\ x_2]$; let $K(\mathbf{x_i},\mathbf{x_j})=(1 + \mathbf{x_i}^T\mathbf{x_j})^2$,**

  **Need to show that $K(\mathbf{x_i},\mathbf{x_j})= \varphi(\mathbf{x_i})^T\varphi(\mathbf{x_j})$:**

  $K(\mathbf{x_i},\mathbf{x_j})=(1 + \mathbf{x_i}^T\mathbf{x_j})^2$,

  $\qquad = 1+ x_{i1}^2 x_{j1}^2 + 2\ x_{i1}x_{j1}\ x_{i2}x_{j2} + x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$

  $\quad = [1\ \ x_{i1}^2\ \ \sqrt{2}\ x_{i1}x_{i2}\ \ x_{i2}^2\ \ \sqrt{2}x_{i1}\ \ \sqrt{2}x_{i2}]^T [1\ \ x_{j1}^2\ \ \sqrt{2}\ x_{j1}x_{j2}\ \ x_{j2}^2\ \ \sqrt{2}x_{j1}\ \ \sqrt{2}x_{j2}]$

  $\quad = \varphi(\mathbf{x_i})^T\varphi(\mathbf{x_j}),\quad$ **where $\varphi(\mathbf{x}) = [1\ \ x_1^2\ \ \sqrt{2}\ x_1x_2\ \ x_2^2\ \ \sqrt{2}x_1\ \ \sqrt{2}x_2]$**

# What Functions are Kernels?

- **For some functions $K(x_i, x_j)$ checking that**

    $$K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j) \text{ can be cumbersome.}$$

- **Mercer's theorem:**

    *Every semi-positive definite symmetric function is a kernel*

- **Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:**

$K=$

| $K(\mathbf{x_1},\mathbf{x_1})$ | $K(\mathbf{x_1},\mathbf{x_2})$ | $K(\mathbf{x_1},\mathbf{x_3})$ | … | $K(\mathbf{x_1},\mathbf{x_N})$ |
|---|---|---|---|---|
| $K(\mathbf{x_2},\mathbf{x_1})$ | $K(\mathbf{x_2},\mathbf{x_2})$ | $K(\mathbf{x_2},\mathbf{x_3})$ | | $K(\mathbf{x_2},\mathbf{x_N})$ |
| … | … | … | … | … |
| $K(\mathbf{x_N},\mathbf{x_1})$ | $K(\mathbf{x_N},\mathbf{x_2})$ | $K(\mathbf{x_N},\mathbf{x_3})$ | … | $K(\mathbf{x_N},\mathbf{x_N})$ |

# Examples of Kernel Functions

- Linear: $K(\mathbf{x_i}, \mathbf{x_j}) = \mathbf{x_i}^\mathsf{T} \mathbf{x_j}$

- Polynomial of power $p$: $K(\mathbf{x_i}, \mathbf{x_j}) = (1 + \mathbf{x_i}^\mathsf{T} \mathbf{x_j})^p$

- Gaussian (radial-basis function network):

$$K(\mathbf{x_i}, \mathbf{x_j}) = \exp(-\frac{\left\| \mathbf{x_i} - \mathbf{x_j} \right\|^2}{2\sigma^2})$$

- Sigmoid: $K(\mathbf{x_i}, \mathbf{x_j}) = \tanh(\beta_0 \mathbf{x_i}^\mathsf{T} \mathbf{x_j} + \beta_1)$

# Non-linear SVMs Mathematically

- **Dual problem formulation:**

> **Find $\alpha_1 \ldots \alpha_N$ such that**
> **$Q(\alpha) = \Sigma \alpha_i - \frac{1}{2} \Sigma \Sigma \alpha_i \alpha_j y_i y_j K(\mathbf{x_i}, \mathbf{x_j})$ is maximized and**
> **(1) $\Sigma \alpha_i y_i = 0$**
> **(2) $\alpha_i \geq 0$ for all $\alpha_i$**

- **The solution is:**

$$f(\mathbf{x}) = \Sigma \alpha_i y_i K(\mathbf{x_i}, \mathbf{x_j}) + b$$

- **Optimization techniques for finding $\alpha_i$'s remain the same!**

# Nonlinear SVM - Overview

- SVM locates a separating hyperplane in the feature space and classify points in that space

- It does not need to represent the space explicitly, simply by defining a kernel function

- The kernel function plays the role of the dot product in the feature space.

# Properties of SVM

- **Flexibility in choosing a similarity function**
- **Sparseness of solution when dealing with large data sets**

  **- only support vectors are used to specify the separating hyperplane**

- **Ability to handle large feature spaces**

  **- complexity does not depend on the dimensionality of the feature space**

- **Overfitting can be controlled by soft margin approach**

- **Nice math property:** **a simple convex optimization problem which is guaranteed to converge to a single global solution**

- **Feature Selection**

# SVM Applications

- **SVM has been used successfully in many real-world problems**
  - **- text (and hypertext) categorization**
  - **- image classification**
  - **- bioinformatics (Protein classification, Cancer classification)**
  - **- hand-written character recognition**

# Weakness of SVM

- **It is sensitive to noise**

  **- A relatively small number of mislabeled examples can dramatically decrease the performance**

- **It only considers two classes**

  **- how to do multi-class classification with SVM?**

  **- Answer:**

  **1) with output arity m, learn m SVM's**

  - **SVM 1 learns "Output==1" vs "Output != 1"**
  - **SVM 2 learns "Output==2" vs "Output != 2"**
  - **:**
  - **SVM m learns "Output==m" vs "Output != m"**

  **2)To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.**

# Back to the Perceptron Example

# The Perceptron

- **The primal algorithm in the feature space**

  $D = \{(x_i, y_i), i = 1, \ldots, m\}$ training data set.

  $\mathbf{x}_i = \phi(x_i) \in \mathcal{K} \subset \mathbb{R}^n$ feature map.

  1., $\mathbf{w} = 0 \in \mathbb{R}^n$

  2., $\forall\ (x_i, y_i),\ i = 1, \ldots, m$, evaluate $sign(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)$

  3., If $x_i$ is misclassified $(sign(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) < 0)$
  $$\text{then } \mathbf{w} := \mathbf{w} + y_i \mathbf{x}_i$$

  4., If no mistakes occur $\Rightarrow$ STOP

# The primal algorithm in the feature space

**Algorithm 1** Perceptron learning algorithm (in primal variables).

**Require:**    A feature mapping $\boldsymbol{\phi} : \mathcal{X} \to \mathcal{K} \subseteq \ell_2^n$

**Ensure:**    A linearly separable training sample $z = ((x_1, y_1), \ldots, (x_m, y_m))$

   $\mathbf{w}_0 = \mathbf{0}; \; t = 0$

   **repeat**

      **for** $j = 1, \ldots, m$ **do**

         **if** $y_j \langle \boldsymbol{\phi}(x_j), \mathbf{w} \rangle \leq 0$ **then**

            $\mathbf{w}_{t+1} = \mathbf{w}_t + y_j \boldsymbol{\phi}(x_j)$

            $t \leftarrow t + 1$

         **end if**

      **end for**

   **until** no mistakes have been made within the **for** loop

   **return** the final weight vector $\mathbf{w}_t$

If $x_j$ is misclassified

# The Perceptron

We start at $\mathbf{w}_0 = 0 \in \mathcal{K} \subset \mathbb{R}^n$

$m=$ num of training examples,
$n = dim(\mathcal{K})$,

$t=$ num of mistakes so far

$$\Rightarrow \mathbf{w}_t = \sum_{i=1}^{m} \alpha_i \phi(x_i) = \sum_{i=1}^{m} \alpha_i \mathbf{x}_i \in \mathbb{R}^n \text{ at time step } t$$

Thus instead of tuning $n$ variables
$$\mathbf{w} = (w_1, \dots, w_n) \text{ (\textbf{Primal variables})}$$
in the large $n$-dimensional feautre space $\mathcal{K}$, it is
enough to learn $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)$ values (**Dual variables**).

# The Perceptron

## The Dual Algorithm in the feature space

$D = \{(x_i, y_i), i = 1, \ldots, m\}$ training data set.

$\mathbf{x}_i = \phi(x_i) \in \mathcal{K} \subset \mathbb{R}^n$ feaure map, $i = 1, \ldots, m$

$t =$ num of mistakes so far

$\Rightarrow \mathbf{w}_t = \sum_{i=1}^{m} \alpha_i \phi(x_i) = \sum_{i=1}^{m} \alpha_i \mathbf{x}_i \in \mathbb{R}^n$ at time step $t$

We update $\boldsymbol{\alpha_t} \in \mathbb{R}^m$ whenever a mistake occurs

1., $\boldsymbol{\alpha}_0 = 0 \in \mathbb{R}^m$

2., $\forall j = 1, \ldots, m$ evaluate

$$y_j \langle \mathbf{x}_j, \mathbf{w}_t \rangle = y_j \langle \mathbf{x}_j, \sum_{i=1}^{m} \alpha_i \mathbf{x}_i \rangle = y_j \sum_{i=1}^{m} \alpha_i \langle \mathbf{x}_j, \mathbf{x}_i \rangle$$

3., If $x_j$ is misclassified ($y_j \langle \mathbf{x}_j, \mathbf{w}_t \rangle < 0$) then update $\boldsymbol{\alpha}_t \in \mathcal{K}$

4., If no mistakes occur $\Rightarrow$ STOP

# The Dual Algorithm in the feature space

**Algorithm 2** Perceptron learning algorithm (in dual variables).

**Require:** A feature mapping $\phi : \mathcal{X} \to \mathcal{K} \subseteq \ell_2^n$

**Ensure:** A linearly separable training sample $z = ((x_1, y_1), \ldots, (x_m, y_m))$

$\alpha = 0$

**repeat**

    **for** $j = 1, \ldots, m$ **do**

        **if** $y_j \sum_{i=1}^{m} \alpha_i \langle \phi(x_i), \phi(x_j) \rangle \leq 0$ **then**    If $x_j$ is misclassified

            $\alpha_j \leftarrow \alpha_j + y_j$

        **end if**

    **end for**

**until** no mistakes have been made within the **for** loop

**return** the vector $\alpha$ of expansion coefficients

# The Dual Algorithm in the feature space

For the classification of a new object $(x, y)$
we have to evaluate

$$y \sum_{i=1}^{m} \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

We don't have to know the actual values of $\mathbf{x} = \phi(x)$!

It is enough to know the inner products

$$\langle \mathbf{x}, \mathbf{x}_i \rangle \quad \forall i = 1, \ldots, m$$

between the object and the training points

# Kernels

**Definition**: **(kernel)**

We are given $\phi : \mathcal{X} \to \mathcal{K} \subset l_2^n$ feautre mapping.

The **kernel** $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is the corresponding inner product function:

$$k(x_i, x_j) \doteq \langle \underbrace{\phi(x_i)}_{\mathbf{x}_i}, \underbrace{\phi(x_j)}_{\mathbf{x}_j} \rangle_{\mathcal{K}} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle_{\mathcal{K}}$$

# Kernels

**Definition**: **(Gram matrix, kernel matrix)**

Gram matrix $G \in \mathbb{R}^{m \times m}$ of kernel $k$ at $\{x_1, \ldots, x_m\}$ :

$$\left. \begin{array}{l} \text{Given a kernel } k : \mathcal{X} \times \mathcal{X} \to \mathbb{R} \\ \text{and a training set } \{x_1, \ldots, x_m\} \end{array} \right\} \Rightarrow G_{ij} \doteq k(x_i, x_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

**Definition**: **(Feature space, kernel space)**

$$\mathcal{K} \doteq span\{\phi(x) \mid x \in \mathcal{X}\} \subset \mathbb{R}^n$$

# Kernel technique

## Definition:

Matrix $G \in \mathbb{R}^{m \times m}$ is positive semidefinite (PSD)
$$\Leftrightarrow G \text{ is symmetric, and } 0 \leq \boldsymbol{\beta}^T G \boldsymbol{\beta} \ \forall \boldsymbol{\beta} \in \mathbb{R}^{m \times m}$$

$$\left. \begin{array}{l} \text{Given a kernel } k : \mathcal{X} \times \mathcal{X} \to \mathbb{R} \\ \text{and a training set } \{x_1, \ldots, x_m\} \end{array} \right\} \Rightarrow G_{ij} \doteq k(x_i, x_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle_{\mathcal{K}}$$

## Lemma:

The Gram matrix is symmetric, PSD matrix.

## Proof:

$$\mathbf{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_m] \in \mathbb{R}^{n \times m} \Rightarrow G = \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{m \times m}$$

$$0 \leq \langle \mathbf{X}\boldsymbol{\beta}, \mathbf{X}\boldsymbol{\beta} \rangle_{\mathcal{K}} = \boldsymbol{\beta}^T G \boldsymbol{\beta}$$

# Kernel technique

We already know that several algorithms use the **kernel values** only (…and NOT the **feature values**)!

**Key idea:**

Choose a nice kernel function $k$
rather than an ugly feautre mapping
$\phi : \mathcal{X} \to \mathbb{R}^n$

# Kernel technique

We have seen so far how to build a kernel $k(\cdot, \cdot)$ from a given feature map $\phi : \mathcal{X} \to \mathbb{R}^n$

Now we want to do the opposite:

A function $k(\cdot, \cdot)$ is kernel $\Leftrightarrow$ there exists a feature space $\mathcal{K}$ and feature map $\phi : \mathcal{X} \to \mathcal{K}$, such that $k(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle_{\mathcal{K}}$

## Let us try to find $\phi$ and $\mathcal{K}$!

# From Vector domain to Functions

- Observe that each vector v = (v[1], v[2], ..., v[n])
  is a mapping from the integers {1,2,..., n} to <

- We can generalize this easily to **INFINITE** domain
  w = (w[1], w[2], ..., w[n], ...)
  where w is mapping from {1,2,...} to <



$$(T_G v)(i) \doteq (Gv)(i) = \underbrace{\sum_{j=1}^{\infty}}_{\int_{\mathcal{X}}} \underbrace{G_{ij}}_{k(i,j)} \underbrace{v_j}_{f(j)}$$

# From Vector domain to Functions

From integers we can further extend to

- < or
- $<^m$
- Strings
- Graphs
- Sets
- Whatever
- …

# Kernels

We don't need the $\mathcal{K} \subset l_2^n$ assumption. It is enough if $\mathcal{K}$ is a complete inner product (Hilbert) space.

## Definition: inner product, Hilbert spaces

$\langle \cdot, \cdot \rangle : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$ is an inner product in vector space $\mathcal{K}$, iff for all vectors $x, y, z \in \mathcal{K}$ and all scalars $a \in \mathbb{R}$:

* Symmetry: $\langle x, y \rangle = \langle y, x \rangle$.

* Linearity in the first argument:
$\langle ax, y \rangle = a\langle x, y \rangle, \ \langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$.

* Positive-definite: $\langle x, x \rangle \geq 0$ with equality only for $x = 0$.

This is more general than the inner product in $\mathbb{R}^n = l_2^n$

**Examples**:
- space of square integrable functions $L_2(\mathcal{X})$,
- space of square summable infinite series $l_2$

# Integral operators, eigenfunctions

## Definition: Eigenvalue, Eigenfunction

- $\lambda$ is the eigenvalue,
- $\Psi \in L_2(\mathcal{X})$ is the eigenfunction

    of integral opreator $(T_k f)(\cdot) = \int\limits_{\mathcal{X}} k(\cdot, x) f(x) dx$

$$\Leftrightarrow \begin{cases} \int\limits_{\mathcal{X}} k(x, \bar{x}) \psi(\bar{x}) d\bar{x} = \lambda \psi(x) \quad \forall x \in \mathcal{X} \\ \\ \|\psi\|_{L_2}^2 \doteq \int\limits_{\mathcal{X}} \psi^2(x) dx = 1 \end{cases}$$

The previous $Gv = \lambda v$ is a special case of this, when $\mathcal{X} = \{x_1, \ldots, x_r\}$ is a finite set.

# Positive (semi) definite operators

**Definition: Positive Definite Operator**

$k(\cdot, \cdot)$ is symmetric kernel,

$$\Rightarrow (T_k f)(\cdot) \doteq \int_{\mathcal{X}} k(\cdot, x) f(x) dx$$

$T_k : L_2(\mathcal{X}) \to L_2(\mathcal{X})$ operator is positive semi definit

$$\Leftrightarrow \int_{\mathcal{X}} \int_{\mathcal{X}} k(\tilde{x}, x) f(x) f(\tilde{x}) dx d\tilde{x} \geq 0 \quad \forall f \in L_2(\mathcal{X})$$

The previous $v^T G v \geq 0$ is a special case of this, when $\mathcal{X} = \{x_1, \ldots, x_r\}$ is a finite set.

# Mercer's theorem

$$(*) \begin{cases} k(\cdot, \cdot) \in L_2(\mathcal{X} \times \mathcal{X}), \\[2mm] k \text{ is symmetric: } k(x, \tilde{x}) = k(\tilde{x}, x) \\[2mm] (T_k f)(\cdot) = \int_{\mathcal{X}} k(\cdot, x) f(x) dx \text{ operator is pos. semi definit} \\[2mm] \psi_i, \ i = 1, 2, \dots \text{ are the eigenfunctions of } T_k \\ \text{with eigenvalues } \lambda_i \end{cases}$$

$$\Rightarrow \begin{cases} (\lambda_1, \lambda_2, \dots) \in l_1, \quad \lambda_i \geq 0 \ \forall i \\[2mm] \psi_i \in L_\infty(\mathcal{X}), \quad \forall i = 1, 2, \dots \\[2mm] k(x, \tilde{x}) = \sum\limits_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(\tilde{x}) \quad \forall x, \tilde{x} \end{cases}$$

2 variables         1 variable

# Mercer's theorem

We like the Mercer's theorem becuase of the **expansion**:

$$k(x, \tilde{x}) = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(\tilde{x}) \quad \forall x, \tilde{x}$$

It shows the **existence of the feature map** $\phi : \mathcal{X} \to \mathcal{K} \subset l_2$

Let $\mathcal{K} \doteq l_2$,
and let $\phi(x) \doteq (\sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x), \dots)^T$

$$\Rightarrow \langle \phi(x), \phi(\tilde{x}) \rangle_{l_2}$$
$$= (\sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x), \dots)^T (\sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x), \dots)$$
$$= \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(\tilde{x}) = k(x, \tilde{x}) \dots \smile$$

$\psi(x) = (\psi_1(x), \psi_2(x), \dots)$ is known as **Mercer map**

# A nicer characterization

The (*) condition in the Mercer's theorem is a bit ugly, but we have a nicer form that characterizes when a function $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a kernel
(i.e. scalar product in some inner product space)

**Theorem:** nicer kernel characterization

$k(\cdot, \cdot)$ is a (Mercer) kernel

$\Leftrightarrow (T_k f)(\cdot)$ is a pos. semi definite operator

$\Leftrightarrow G = (k(x_i, x_j))_{ij}^r \in \mathbb{R}^{r \times r}$ Gram matrix is pos. semi definite $\forall r,\ \forall (x_1, \ldots, x_r) \in \mathcal{X}^r$

# Vapnik-Chervonenkis dimension

$$\text{TESTERR}(\alpha) = E\left[\frac{1}{2}|y - f(x, \alpha)|\right] \qquad \text{TRAINERR}(\alpha) = \frac{1}{R}\sum_{k=1}^{R}\frac{1}{2}|y_k - f(x_k, \alpha)|$$

- Given some machine **f**, let *h* be its VC dimension.
- *h* is a measure of **f**'s power (*h* does not depend on the choice of training set)
- Vapnik showed that with probability 1-η

$$\text{TESTERR}(\alpha) \leq \text{TRAINERR}(\alpha) + \sqrt{\frac{h(\log(2R/h)+1) - \log(\eta/4)}{R}}$$

This gives us a way to estimate the error on future data based only on the training error and the VC-dimension of **f**

# Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?



- How about… mapping data to a higher-dimensional space:

# The "Kernel Trick"

- **To produce linear separability in Higher Dimension, the linear classifier relies on dot product between vectors $K(\mathbf{x_i},\mathbf{x_j})=\mathbf{x_i^T}\mathbf{x_j}$**

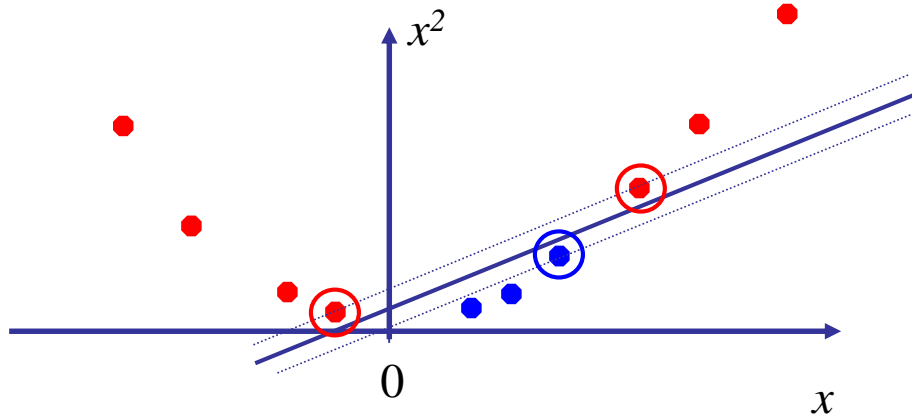- **If every data point is mapped into high-dimensional space via some transformation $\Phi$: $\mathbf{x} \rightarrow \varphi(\mathbf{x})$, the dot product becomes:**

$$K(\mathbf{x_i},\mathbf{x_j})= \varphi(\mathbf{x_i})^{T}\varphi(\mathbf{x_j})$$

- **A *kernel function* is some function that corresponds to an inner product in some expanded feature space.**

- **Example:**

  **2-dimensional vectors $\mathbf{x}=[x_1\ \ x_2]$; let $K(\mathbf{x_i},\mathbf{x_j})=(1 + \mathbf{x_i^T}\mathbf{x_j})^2$,**

  **Need to show that $K(\mathbf{x_i},\mathbf{x_j})= \varphi(\mathbf{x_i})^{T}\varphi(\mathbf{x_j})$:**

  $K(\mathbf{x_i},\mathbf{x_j})=(1 + \mathbf{x_i^T}\mathbf{x_j})^2$,

  $= 1+ x_{i1}^2 x_{j1}^2 + 2\ x_{i1}x_{j1}\ x_{i2}x_{j2}+ x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$

  $= [1\ \ x_{i1}^2\ \sqrt{2}\ x_{i1}x_{i2}\ \ x_{i2}^2\ \sqrt{2}x_{i1}\ \sqrt{2}x_{i2}]^T [1\ \ x_{j1}^2\ \sqrt{2}\ x_{j1}x_{j2}\ \ x_{j2}^2\ \sqrt{2}x_{j1}\ \sqrt{2}x_{j2}]$

  $= \varphi(\mathbf{x_i})^{T}\varphi(\mathbf{x_j})$,   **where $\varphi(\mathbf{x}) = [1\ \ x_1^2\ \sqrt{2}\ x_1x_2\ \ x_2^2\ \sqrt{2}x_1\ \sqrt{2}x_2]$**

# Examples of Kernel Functions

- Linear: $K(\mathbf{x_i}, \mathbf{x_j}) = \mathbf{x_i}^\mathsf{T} \mathbf{x_j}$

- Polynomial of power $p$: $K(\mathbf{x_i}, \mathbf{x_j}) = (1 + \mathbf{x_i}^\mathsf{T} \mathbf{x_j})^p$

- Gaussian (radial-basis function network):

$$K(\mathbf{x_i}, \mathbf{x_j}) = \exp(-\frac{\left\| \mathbf{x_i} - \mathbf{x_j} \right\|^2}{2\sigma^2})$$

- Sigmoid: $K(\mathbf{x_i}, \mathbf{x_j}) = \tanh(\beta_0 \mathbf{x_i}^\mathsf{T} \mathbf{x_j} + \beta_1)$

# Mercer's theorem

$(*)$
$$
\begin{cases}
k(\cdot, \cdot) \in L_2(\mathcal{X} \times \mathcal{X}), \\[2ex]
k \text{ is symmetric: } k(x, \tilde{x}) = k(\tilde{x}, x) \\[2ex]
(T_k f)(\cdot) = \int_{\mathcal{X}} k(\cdot, x) f(x) dx \text{ operator is pos. semi definit} \\[2ex]
\psi_i, \ i = 1, 2, \ldots \text{ are the eigenfunctions of } T_k \\
\text{with eigenvalues } \lambda_i
\end{cases}
$$

$\Rightarrow$
$$
\begin{cases}
(\lambda_1, \lambda_2, \ldots) \in l_1, \quad \lambda_i \geq 0 \ \forall i \\[2ex]
\psi_i \in L_\infty(\mathcal{X}), \quad \forall i = 1, 2, \ldots \\[2ex]
k(x, \tilde{x}) = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(\tilde{x}) \quad \forall x, \tilde{x}
\end{cases}
$$

2 variables          1 variable

# Reproducing Kernel Hilbert Spaces

For a given kernel $k(\cdot, \cdot)$ we already know how to define feature space $\mathcal{K}$, and $\phi : \mathcal{X} \to \mathcal{K}$ feature map (Mercer map):

$$\mathcal{K} = l_2, \text{ and } \phi(x) \doteq (\sqrt{\lambda_1}\psi_1(x), \sqrt{\lambda_2}\psi_2(x), \ldots)^T$$

## Now, we show another way using RKHS

$k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ given kernel $\Rightarrow \mathcal{F}_0 \doteq \{k(x, \cdot) | x \in \mathcal{X}\}$ function space

We will add inner product to $\mathcal{F}_0$ function space
$\Rightarrow$ Pre-Hilbert space

Completing (closing) a pre-Hilbert space $\rightarrow$ Hilbert space
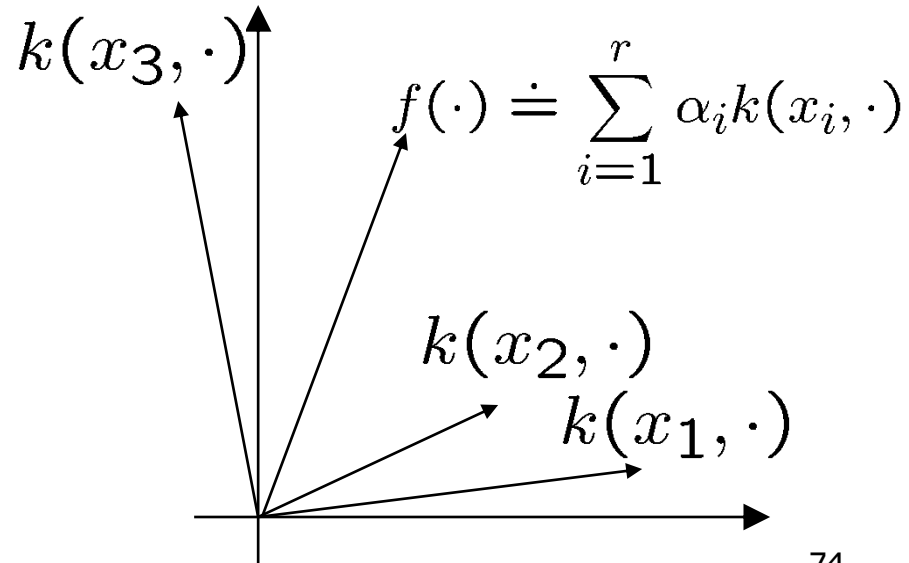
# Reproducing Kernel Hilbert Spaces

$k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ given kernel $\Rightarrow \mathcal{F}_0 \doteq \{k(x, \cdot) | x \in \mathcal{X}\}$ function space

$$(x_1, \ldots, x_r) \text{ given} \Rightarrow f(\cdot) \doteq \sum_{i=1}^{r} \alpha_i k(x_i, \cdot) \ \in \mathcal{F}_0$$

$$(\tilde{x}_1, \ldots, \tilde{x}_s) \text{ given} \Rightarrow g(\cdot) \doteq \sum_{j=1}^{s} \beta_j k(\tilde{x}_j, \cdot) \ \in \mathcal{F}_0$$

**The inner product:**

$$\langle f, g \rangle_{\mathcal{F}_0} \doteq \sum_{i=1}^{r} \sum_{j=1}^{s} \alpha_i \beta_j k(x_i, \tilde{x}_j)$$

$$= \sum_{i=1}^{r} \alpha_i g(x_i)$$

$$= \sum_{j=1}^{s} \beta_j f(\tilde{x}_j) \quad (*)$$

$k(x_3, \cdot)$

$f(\cdot) \doteq \sum_{i=1}^{r} \alpha_i k(x_i, \cdot)$

$k(x_2, \cdot)$

$k(x_1, \cdot)$

# Reproducing Kernel Hilbert Spaces

**Note:**

While for calculating $\langle f, g \rangle_{\mathcal{F}_0}$ we use their representations: $\alpha \in \mathbb{R}^r, \beta \in \mathbb{R}^s, \{x_i\}_{i=1}^r, \{\tilde{x}_j\}_{j=1}^s$ the $\langle f, g \rangle_{\mathcal{F}_0}$ is independent of the representation of $f, g$

**Proof:**

If we change $\alpha \in \mathbb{R}^r$ or $x_i \Rightarrow \langle f, g \rangle_{\mathcal{F}_0}$ doesn't change (because of (*)) The same for $\beta \in \mathbb{R}^s$

$$\langle f, g \rangle_{\mathcal{F}_0} = \sum_{i_1}^{r} \alpha_i f(x_i) = \sum_{j=1}^{s} \beta_j f(\tilde{x}_j) \quad (*)$$

# Reproducing Kernel Hilbert Spaces

**Lemma:**

$\langle f, g \rangle$ is an inner product of $\mathcal{F}_0$

$\qquad \Rightarrow \mathcal{F}_0$ is pre-Hilbert space

$\qquad \mathcal{F} \doteq close(\mathcal{F}_0)$ is a Hilbert space

- **Pre-Hilbert** space:
    Like the Euclidean space with *rational* scalars only

- **Hilbert space**:
    Like the Euclidean space with *real* scalars

**Proof:**

1., $\langle f, g \rangle_{\mathcal{F}_0} = \langle g, f \rangle_{\mathcal{F}_0}$

2., $\langle cf + dg, h \rangle_{\mathcal{F}_0} = c\langle f, h \rangle_{\mathcal{F}_0} + d\langle g, h \rangle_{\mathcal{F}_0}, \ \forall c, d \in \mathbb{R}, \ \forall f, g, h \in \mathcal{F}_0$

3., $\langle f, f \rangle_{\mathcal{F}_0} \geq 0$

4., $\langle f, f \rangle_{\mathcal{F}_0} = 0 \Leftrightarrow f = 0$

# Reproducing Kernel Hilbert Spaces

**Lemma: (Reproducing property)**

$$\langle f, k(x, \cdot)\rangle_{\mathcal{F}} = f(x)$$

**Proof**: definition of $\langle f, g\rangle_{\mathcal{F}}$

**Lemma:** The constructed features match to *k*

Huhh...

$$\langle \underbrace{k(x_i, \cdot)}_{\phi(x_i)}, \underbrace{k(x_j, \cdot)}_{\phi(x_j)}\rangle_{\mathcal{F}} = k(x_i, x_j)$$

**Proof**: reproducing property

# Reproducing Kernel Hilbert Spaces

**Proof of property 4.,:**

$$0 \ \leq (f(x))^2 = \langle f, k(x, \cdot) \rangle_{\mathcal{F}}^2, \ \forall x$$

rep. property

$$\langle f, k(x, \cdot) \rangle_{\mathcal{F}}^2 \leq \langle f, f \rangle_{\mathcal{F}} \langle k(x, \cdot), k(x, \cdot) \rangle_{\mathcal{F}} \quad \forall x$$

we need only that <0,0>=0!

$$\text{Hence, if } \langle f, f \rangle_{\mathcal{F}} = 0 \Rightarrow (f(x))^2 = 0, \ \forall x \in \mathcal{X}$$

$$\Rightarrow f(x) = 0, \ \forall x \in \mathcal{X}$$

$$\Rightarrow f = 0$$

# The Representer Theorem

In the perceptron problem we could use the dual algorithm, because we had this representation:

$$f(x) \doteq \langle \phi(x), \mathbf{w} \rangle_{\mathcal{K}} = \sum_{i=1}^{m} \alpha_i k(x, x_i)$$

and thus we had to update $\alpha_1, \ldots, \alpha_m$ only, and not $\mathbf{w} \in \mathcal{K}$!

The **Representer theorem** provides us a big class of problems, where the solution can be represented by

$$f(\cdot) = \sum_{i=1}^{m} \alpha_i k(x_i, \cdot), \quad \alpha \in \mathbb{R}^m$$

# The Representer Theorem

**Theorem:**

$$k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \to \mathbb{R}, \text{Mercer kernel on } \mathcal{X}$$

$$z = (x_1, y_1), \ldots, (x_m, y_m) \in (\mathcal{X} \times \mathcal{Y})^m \text{ training sample}$$

$$g_{emp} : (\mathcal{X} \times \mathcal{Y} \times \mathbb{R})^m \to \mathbb{R} \cup \{\infty\} \quad \Bigg\} \Rightarrow$$

$$g_{reg} : \mathbb{R} \to [0, \infty) \text{ strictly increasing function}$$

$$\mathcal{F} : \text{ RKHS induced by } k(\cdot, \cdot)$$

$$\Rightarrow f^* = \arg\min_{f \in \mathcal{F}} R_{reg}[f, z]$$

$$\doteq \arg\min_{f \in \mathcal{F}} \underbrace{g_{emp}[(x_i, y_i, f(x_i))_{i \in \{1\ldots m\}}]}_{} + \underbrace{g_{reg}(\|f\|)}_{}$$

1st term, empirical loss    2nd term, regularization

admits the following representation:

$$f^*(\cdot) = \sum_{i=1}^{m} \alpha_i k(x_i, \cdot), \quad \alpha = (\alpha_1, \ldots, \alpha_m) \in \mathbb{R}^m$$

# The Representer Theorem

**Message:**
  *Optimizing in general function classes is difficult, but in RKHS it is only finite! (m) dimensional problem*

**Proof of Representer Theorem:**

$$\phi(x) \doteq k(x, \cdot) = \phi(x)(\cdot)$$

$x_1, \ldots, x_m$ training samples are given

$$f \in \mathcal{F} \Rightarrow f(\cdot) = \sum_{i=1}^{m} \alpha_i \phi(x_i)(\cdot) + v(\cdot)$$

  where $\mathcal{F} \ni v \perp span\{\phi(x_1), \ldots, \phi(x_m)\}$,

  thus $\langle v, \phi(x_i) \rangle_{\mathcal{F}} = 0 \quad \forall i = 1, \ldots, m$

# Proof of the Representer Theorem

**Proof of Representer Theorem**

$$f^* = \arg\min_{f \in \mathcal{F}} R_{reg}[f, z] \doteq \arg\min_{f \in \mathcal{F}} \underbrace{g_{emp}[(x_i, y_i, f(x_i))_{i \in \{1...m\}}]}_{} + \underbrace{g_{reg}(\|f\|)}_{}$$

1st term, empirical loss      2nd term, regularization

$$\Rightarrow f(x_j) = \langle f, \underbrace{k(x_j, \cdot)}_{\phi(x_j)} \rangle_{\mathcal{F}} = \langle \sum_{i=1}^{m} \alpha_i \phi(x_i) + v, \phi(x_j) \rangle_{\mathcal{F}}$$

$$= \sum_{i=1}^{m} \alpha_i \langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{F}} = \sum_{i=1}^{m} \alpha_i k(x_i, x_j)$$

$\Rightarrow f(x_j)$ depends only on $\alpha_1, \ldots, \alpha_m$, but independent from $v$!

$\Rightarrow 1^{st}$ term depends only on $\alpha_1, \ldots, \alpha_m$, but not on $v$

# Proof of the Representer Theorem

$$f^* = \arg\min_{f \in \mathcal{F}} R_{reg}[f, z] \doteq \arg\min_{f \in \mathcal{F}} \underbrace{g_{emp}[(x_i, y_i, f(x_i))_{i \in \{1...m\}}]}_{} + \underbrace{g_{reg}(\|f\|)}_{}$$

1st term, empirical loss        2nd term, regularization

Let us examine the $2^{nd}$ term.

$$g_{reg}(\|f\|) = g_{reg}(\|\sum_{i=1}^{m} \alpha_i \phi(x_i) + v\|)$$

$$= g_{reg}(\sqrt{\|\sum_{i=1}^{m} \alpha_i \phi(x_i)\|_{\mathcal{F}}^2 + \|v\|_{\mathcal{F}}^2}$$

since $\mathcal{F} \ni v \perp span\{\phi(x_1), \ldots, \phi(x_m)\}$

$$\geq g_{reg}(\|\sum_{i=1}^{m} \alpha_i \phi(x_i)\|_{\mathcal{F}})$$

with equality only if $v = 0$!

$\Rightarrow$ any minimizer $f^*$ must have $v = 0$

$$\Rightarrow f^*(\cdot) = \sum_{i=1}^{m} \alpha_i k(x_i, \cdot)$$