# ISTA 421 + INFO 521 Introduction to Machine Learning

**Lecture 23:**
**Support Vector Machines - II and The Kernel Trick**

**Clay Morrison**

claytonm@email.arizona.edu

Harvill 437A

Phone 621-6609

15 November 2017

SISTA  1

---

## Maximum Margin

$\alpha$

$x \longrightarrow$ $f$ $\longrightarrow t$

$t_{\text{new}} = \text{sign}(\mathbf{w}^\mathsf{T}\mathbf{x}_{\text{new}} + b)$

· denotes +1
○ denotes -1

Support Vectors are those datapoints that the margin pushes up against

The maximum margin linear classifier is the linear classifier with the, um, maximum margin.
This is the simplest kind of SVM (Called an LSVM)

Linear SVM

SISTA  2

# Defining the Margin

$$\mathbf{a} \cdot \frac{\mathbf{b}}{\|\mathbf{b}\|} = \|\mathbf{a}\| \cos\theta \quad \|x\| := \sqrt{x \cdot x}.$$

Our "*decision function*", which we will also refer to as $D(x)$:

$$t_{\text{new}} = \text{sign } (\mathbf{w}^\mathsf{T} \mathbf{x}_{\text{new}} + b)$$

**Another complication:**

We note that the argument in $D(x)$ is invariant under a rescaling: $\mathbf{w} \to \lambda\mathbf{w},\ b \to \lambda b.$

We will implicitly fix a scale with:

$$\mathbf{w} \cdot \mathbf{x_1} + b = 1$$
$$\mathbf{w} \cdot \mathbf{x_2} + b = -1$$

for the support vectors (*canonical hyperplanes*).

$\mathbf{w}$ is in the direction **perpendicular** to the boundary.
Normalize it to get the "unit vector": $\dfrac{\mathbf{w}}{\|\mathbf{w}\|}$

Combine both **constraints** by subtracting one from the other, to get the following:

$$\mathbf{w} \cdot (\mathbf{x_1} - \mathbf{x_2}) = 2$$

The margin (the length of the distance between the support vector canonical hyperplanes) will be given by the ***projection*** of the vector $(x_1 - x_2)$ onto the normal vector to the hyperplane!

The projection is accomplished by taking the inner product of these two quantities   ◆SISTA▶ 3

---



# Defining the Margin

$$\mathbf{a} \cdot \frac{\mathbf{b}}{\|\mathbf{b}\|} = \|\mathbf{a}\| \cos\theta \quad \|x\| := \sqrt{x \cdot x}.$$

Our "*decision function*", which we will also refer to as $D(x)$:

$$t_{\text{new}} = \text{sign } (\mathbf{w}^\mathsf{T} \mathbf{x}_{\text{new}} + b)$$

**Another complication:**

We note that the argument in $D(x)$ is invariant under a rescaling: $\mathbf{w} \to \lambda\mathbf{w},\ b \to \lambda b.$

We will implicitly fix a scale with:

$$\mathbf{w} \cdot \mathbf{x_1} + b = 1$$
$$\mathbf{w} \cdot \mathbf{x_2} + b = -1$$

for the support vectors (*canonical hyperplanes*).

$\mathbf{w}$ is in the direction **perpendicular** to the boundary.
Normalize it to get the "unit vector": $\dfrac{\mathbf{w}}{\|\mathbf{w}\|}$

Combine both **constraints** by subtracting one from the other, to get the following:
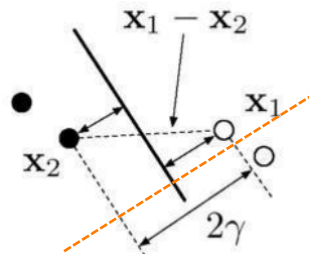
$$\mathbf{w} \cdot (\mathbf{x_1} - \mathbf{x_2}) = 2$$

The margin (the length of the distance between the support vector canonical hyperplanes) will be given by the ***projection*** of the vector $(x_1 - x_2)$ onto the normal vector to the hyperplane!

The projection is accomplished by taking the inner product of these two quantities   ◆SISTA▶ 4

## Defining the Margin

$$\mathbf{a} \cdot \frac{\mathbf{b}}{\|\mathbf{b}\|} = \|\mathbf{a}\| \cos\theta \quad \|\boldsymbol{x}\| := \sqrt{\boldsymbol{x} \cdot \boldsymbol{x}}.$$

$$2\gamma = \frac{1}{\|\mathbf{w}\|} \mathbf{w}^\mathsf{T}(\mathbf{x}_1 - \mathbf{x}_2)$$
$$= \frac{1}{\|\mathbf{w}\|}(\mathbf{w}^\mathsf{T}\mathbf{x}_1 - \mathbf{w}^\mathsf{T}\mathbf{x}_2)$$
$$= \frac{1}{\|\mathbf{w}\|}(\mathbf{w}^\mathsf{T}\mathbf{x}_1 + b - \mathbf{w}^\mathsf{T}\mathbf{x}_2 - b)$$
$$= \frac{1}{\|\mathbf{w}\|}(1 + 1)$$
$$\gamma = \frac{1}{\|\mathbf{w}\|}.$$

**w** is in the direction **perpendicular** to the boundary.
Normalize it to get the "unit vector": $\dfrac{\mathbf{w}}{\|\mathbf{w}\|}$

**Constraints**:
$$\mathbf{w} \cdot \mathbf{x_1} + b = 1$$
$$\mathbf{w} \cdot \mathbf{x_2} + b = -1$$
$$\mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = 2$$
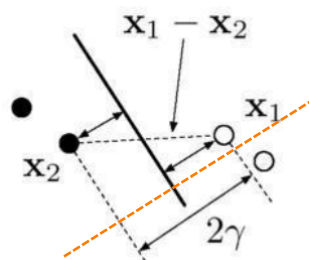
The margin (the length of the distance between the support vector canonical hyperplanes) will be given by the **projection** of the vector $(x_1 - x_2)$ onto the normal vector to the hyperplane!
The projection is accomplished by taking the inner product of these two quantities

SISTA 5

## Maximizing the Margin

$$\gamma = \frac{1}{\|\mathbf{w}\|}$$

Recall, we set a scale for the closest points to the margin:
$$\mathbf{w} \cdot \mathbf{x_1} + b = 1 \quad \text{For the s.v. class 1}$$
$$\mathbf{w} \cdot \mathbf{x_2} + b = -1 \quad \text{For the s.v. class 2}$$

Therefore, **w** must be chosen such that:
$$\mathbf{w} \cdot \mathbf{x_n} + b \geq 1 \quad \text{for all } \mathbf{x_n} \text{ in class 1}$$
$$\mathbf{w} \cdot \mathbf{x_n} + b \leq -1 \quad \text{for all } \mathbf{x_n} \text{ in class 2}$$

Combining both constraints is easy:
$$t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad \text{For all } \mathbf{x}_n$$

There are a total of N constraints
Easier to *minimize* $\dfrac{1}{2}\|\mathbf{w}\|^2$

SISTA 6

# Constrained Optimization with
## Lagrange (KKT) Multipliers

- Find values of a set of parameters that maximize (or minimize) an objective function, but also satisfy some constraints.
- Create new objective function that includes the original plus an additional term for each constraint.

For example, minimize *f(x)* subject to the constraint *g(w)* ≤ *a*

$$\underset{w}{\text{argmin}} \quad f(w)$$

$$\text{subject to} \quad g(w) \leq a$$

Add Lagrange term of the form *λ(a - g(w))* and optimize for *w* and *λ*

$$\underset{w,\lambda}{\text{argmin}} \quad f(w) - \lambda(a - g(w))$$

Bug in book!
(p. 188, comment 5.1)

$$\text{subject to} \quad \lambda > 0$$

(Strictly speaking, Lagrange multipliers are just for equality constraints, whereas constrained optimization problems with inequalities involve KKT (Karush-Kuhn-Tucker) conditions; so you may see the lambda above referred to as a "KKT" multiplier…)

7

# Constrained Optimization with
## Lagrange (KKT) Multipliers



(Shamelessly cribbed from Wikipedia…)

SISTA ▶ 8

# Maximizing the Margin $\gamma = \dfrac{1}{\|\mathbf{w}\|}$

$$\frac{1}{2}\|\mathbf{w}\|^2 \qquad t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$$

Maximizing the margin, γ, becomes a ***constrained optimization*** problem, namely, to minimize the following:

$$\underset{\mathbf{w}}{\mathrm{argmin}} \quad \frac{1}{2}\|\mathbf{w}\|^2$$
$$\text{subject to} \quad t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1, \text{ for all } n$$

We can incorporate the inequalities into the minimization by introducing ***Lagrange multipliers***, resulting in the following:

note: $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{w}$

$$\underset{\mathbf{w},\alpha}{\mathrm{argmin}} \quad \frac{1}{2}\mathbf{w}^\top \mathbf{w} - \sum_{n=1}^{N} \alpha_n(t_n(\mathbf{w}^\top \mathbf{x}_n + b) - 1)$$
$$\text{subject to} \quad \alpha_n \geq 0, \text{ for all } n,$$

Recall how we maximize/minimize!

$$\frac{\partial}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n$$
$$\frac{\partial}{\partial b} = -\sum_{n=1}^{N} \alpha_n t_n.$$

Set to 0!

$$\mathbf{w} = \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n$$
$$\sum_{n=1}^{N} \alpha_n t_n = 0$$

These two identities must be satisfied at the optimum

---

# Maximizing the Margin $\gamma = \dfrac{1}{\|\mathbf{w}\|}$

$$\frac{1}{2}\|\mathbf{w}\|^2 \qquad t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$$

$$\underset{\mathbf{w},\alpha}{\mathrm{argmin}} \quad \frac{1}{2}\mathbf{w}^\top \mathbf{w} - \sum_{n=1}^{N} \alpha_n(t_n(\mathbf{w}^\top \mathbf{x}_n + b) - 1) \qquad \mathbf{w} = \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n$$
$$\text{subject to} \quad \alpha_n \geq 0, \text{ for all } n, \qquad\qquad\qquad \sum_{n=1}^{N} \alpha_n t_n = 0$$

Plug constraint for **w** back into the objective function, to get:

$$\frac{1}{2}\mathbf{w}^\top \mathbf{w} - \sum_{n=1}^{N} \alpha_n(t_n(\mathbf{w}^\top \mathbf{x}_n + b) - 1)$$

$$= \frac{1}{2}\left(\sum_{m=1}^{N} \alpha_m t_m \mathbf{x}_m^\top\right)\left(\sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n\right) - \sum_{n=1}^{N} \alpha_n\left(t_n\left(\sum_{m=1}^{N} \alpha_m t_m \mathbf{x}_m^\top \mathbf{x}_n + b\right) - 1\right)$$

$$= \frac{1}{2}\sum_{n,m=1}^{N} \alpha_m \alpha_n t_m t_n \mathbf{x}_m^\top \mathbf{x}_n - \sum_{n,m=1}^{N} \alpha_m \alpha_n t_m t_n \mathbf{x}_m^\top \mathbf{x}_n - \sum_{n=1}^{N} \alpha_n t_n b + \sum_{n=1}^{N} \alpha_n$$

$$= \sum_{n=1}^{N} \alpha_n - \frac{1}{2}\sum_{n,m=1}^{N} \alpha_m \alpha_n t_m t_n \mathbf{x}_m^\top \mathbf{x}_n$$

This goes away by this constraint

This is the **dual** optimization problem (we have eliminated **w**!)
It is a ***quadratic optimization*** problem due to the $\alpha_m \alpha_n$ term (matlab: quadprog)

SISTA 10

# Making Predictions

Our final constraint problem:

$$\sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n,m=1}^{N} \alpha_m \alpha_n t_m t_n \mathbf{x}_m^{\top} \mathbf{x}_n$$

Subject to: $\alpha_n \geq 0, \ \sum_{n=1}^{N} \alpha_n t_n = 0.$

Give it to a quadratic programming solver!

To predict, we need our decision function D(x)

$$t_{\text{new}} = \text{sign } (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{\text{new}} + b)$$

But we just optimized for α's!

Recall: $\mathbf{w} = \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n$

So rewrite the decision function as:

$$t_{\text{new}} = \text{sign} \left( \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n^{\top} \mathbf{x}_{\text{new}} + b \right)$$

To find $b$, we will use the fact that for the closest points, $t_n \left( \mathbf{w}^{\mathsf{T}} \mathbf{x}_n + b \right) = 1$

$$b = t_n - \sum_{m=1}^{N} \alpha_m t_m \mathbf{x}_m^{\top} \mathbf{x}_n$$

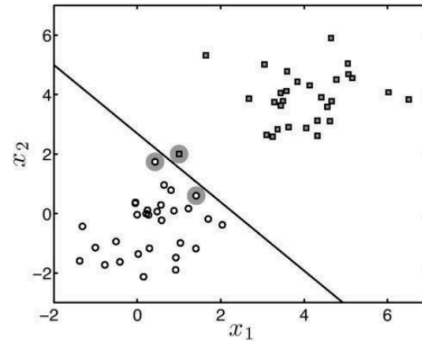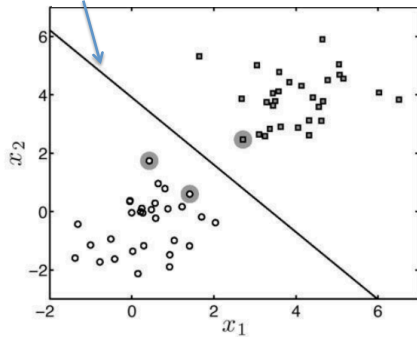(note that $t_n = 1/t_n$ in this case because $t_n = \{ 1, -1 \}$)

SISTA 11

---

# Hard Margin SMV

$$t_{new} = \text{sign} \left( \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n^{\top} \mathbf{x}_{new} + t_n - \sum_{m=1}^{N} \alpha_m t_m \mathbf{x}_m^{\top} \mathbf{x}_n \right)$$

After optimizing for all $\alpha_n$'s, the only α's that are **non-zero** are the **support vectors**!

$(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b = 0, \ \text{where} \ \mathbf{w} = \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n)$



SISTA 12

# Soft Margin SVM

- It allow points to lie on the wrong side of the boundary, need to "slacken" the constraints

$$t_n(\mathbf{w}^\mathsf{T}\mathbf{x}_n + b) \geq 1 - \xi_n \quad \text{where} \quad \xi_n \geq 0$$

- If $0 \leq \xi_n \leq 1$
  - Then the point lies on the correct side of the boundary, but within the boundary margin
- If $\xi_n > 1$
  - Then the point lies on the "wrong" side of the boundary

SISTA 13

# Soft Margin SVM

- It allow points to lie on the wrong side of the boundary, need to "slacken" the constraints

$$t_n(\mathbf{w}^\mathsf{T}\mathbf{x}_n + b) \geq 1 - \xi_n \quad \text{where} \quad \xi_n \geq 0$$

- The optimization task becomes: Recall, the original constrained optimization was:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \frac{1}{2}\mathbf{w}^\mathsf{T}\mathbf{w} + C\sum_{n=1}^{N}\xi_n$$

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \frac{1}{2}\|\mathbf{w}\|^2$$
$$\text{subject to} \quad t_n(\mathbf{w}^\mathsf{T}\mathbf{x}_n + b) \geq 1, \text{ for all } n$$

$$\text{subject to } \xi_n \geq 0 \text{ and } t_n(\mathbf{w}^\mathsf{T}\mathbf{x}_n + b) \geq 1 - \xi_n \text{ for all } n$$

- *C* controls to what extent we are willing to allow points to sit within the margin itself or on the wrong side of the decision boundary

SISTA 14

# Soft Margin SVM

- It allow points to lie on the wrong side of the boundary, need to "slacken" the constraints

$$t_n(\mathbf{w}^\mathsf{T}\mathbf{x}_n + b) \geq 1 - \xi_n \quad \text{where} \quad \xi_n \geq 0$$

- It turns out that incorporating the new constraint $C$ does not change the overall optimization much!: Recall: $\sum_{n=1}^{N}\alpha_n - \frac{1}{2}\sum_{n,m=1}^{N}\alpha_m\alpha_n t_m t_n \mathbf{x}_m^\mathsf{T}\mathbf{x}_n \quad \alpha_n \geq 0, \ \sum_{n=1}^{N}\alpha_n t_n = 0$

$$\underset{\mathbf{w}}{\operatorname{argmax}} \ \sum_{n=1}^{N}\alpha_n - \frac{1}{2}\sum_{n,m=1}^{N}\alpha_n\alpha_m t_n t_m \mathbf{x}_n^\mathsf{T}\mathbf{x}_m$$
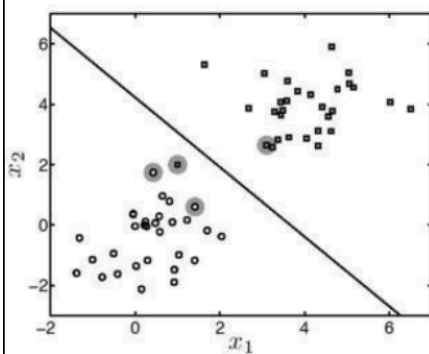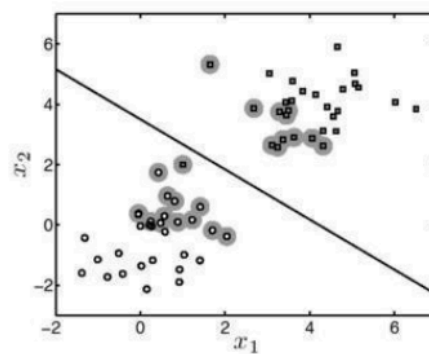
Just adds an upper bound on the influence any training point can have

$$\text{subject to} \ \sum_{n=1}^{N}\alpha_n t_n = 0 \ \text{and} \ 0 \leq \alpha_n \leq C, \ \text{for all } n.$$

SISTA 15

# Decision Boundary & Support Vectors for different $C$



(a) $C = 1$

(b) $C = 0.01$

Larger C = lower bias (fewer support vectors)
Smaller C = higher bias (more support vectors)

$$\underset{\mathbf{w}}{\operatorname{argmax}} \ \sum_{n=1}^{N}\alpha_n - \frac{1}{2}\sum_{n,m=1}^{N}\alpha_n\alpha_m t_n t_m \mathbf{x}_n^\mathsf{T}\mathbf{x}_m$$

$$\text{subject to} \ \sum_{n=1}^{N}\alpha_n t_n = 0 \ \text{and} \ 0 \leq \alpha_n \leq C, \ \text{for all } n.$$

# Kernels: Transforming Data

- So far, the boundary has been linear
- **Recall** in our treatment of Linear Regression, to get non-linear boundaries we added terms to **x** and extended **w**     (… the weight space view)
  - In that case, **we explicitly projected the data**

$$\mathbf{x} \xrightarrow{\phi} \phi(\mathbf{x}) \qquad \text{e.g., } \phi(x) = 1 + x + x^2$$

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

SISTA ⟩ 17

# Kernels: Transforming Data

- Here we take a different approach: the model remains the same (linear decision boundary) but **the data are *implicitly transformed into a new space***.
- **Kernel functions** are *similarity* functions.
  - They compute the similarity of pairs of input points **x** and **z**, <u>as if</u> they had been projected into some space $\nu$ by function $\phi$ (i.e., they compute the similarity between $\phi(\mathbf{x})$ and $\phi(\mathbf{z})$ ), **but without actually doing the projection!**

$$\phi : \mathcal{X} \to \nu \qquad\qquad k(\mathbf{x}, \mathbf{z}) = \underbrace{\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle}_{\text{inner product}}{}_{\nu}$$

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

The word "kernel" is used in mathematics to denote a weighting function for a weighted sum or integral.

Similarity, here, is in terms of an inner product (on $\nu$). An explicit representation for $\phi$ is not necessary **as long as $\nu$ is an inner product space**.

Formally, the real-valued function $k(\mathbf{x},\mathbf{z})$ must satisfy Mercer's condition:

$$\iint g(\mathbf{x})k(\mathbf{x}, \mathbf{z})g(\mathbf{z}) \, d\mathbf{x} \, d\mathbf{z} \geq 0 \quad \text{for all square integrable fns } g(\mathbf{x})$$

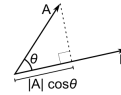$$\int_{-\infty}^{\infty} |g(\mathbf{x})|^2 \, d\mathbf{x} \neq \infty$$

SISTA ⟩ 18

# Kernels: Transforming Data

- The most familiar (from this class so far) is the dot product, where the mapping function $\phi$ is the *identity* function:

$$\phi(\mathbf{x}) = \mathbf{x}$$
$$k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z} = x_1 z_1 + x_2 z_2 + ... + x_n z_n$$

SISTA ▶ 19

# Kernels: Transforming Data

- Another example! $k(\mathbf{x}, \mathbf{z}) = \left(\mathbf{x}^\top \mathbf{z}\right)^2$

  **Find the explicit projection of the vectors**

input space:
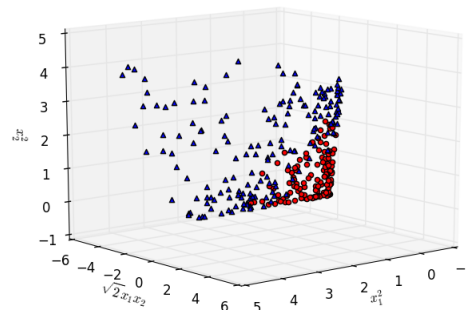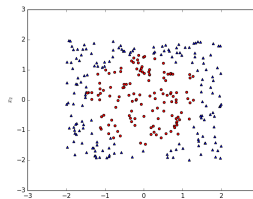$$\mathbf{x} = (x_1, x_2)^\top$$

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

$$\phi(\mathbf{x}) = \left(x_1^2, \sqrt{2} x_1 x_2, x_2^2\right)^\top$$

# The Kernel Trick

- SVMs are one of a **class** of algorithms that don't actually need to perform that actual transformation!
- The terms representing the data only appear within inner (i.e., dot) products: $\mathbf{x}_n^\top \mathbf{x}_m, \ \mathbf{x}_n^\top \mathbf{x}_{new}$

$$\sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n,m=1}^{N} \alpha_m \alpha_n t_m t_n \mathbf{x}_m^\top \mathbf{x}_n \quad \longleftarrow \quad \text{objective}$$

$$\text{prediction} \longrightarrow \quad t_{new} = \text{sign}\left( \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n^\top \mathbf{x}_{new} + t_n - \sum_{m=1}^{N} \alpha_m t_m \mathbf{x}_m^\top \mathbf{x}_n \right)$$

- We never see **x** on its own.
- Were we to do a transformation, we would need to calculate inner products in the new space
- BUT, if we can find a function $k(\mathbf{x}_n, \ \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$ then we can use *k* directly without performing the transform.

SISTA ⟩ 21

# The Kernelized Soft Margin SVM

Original soft margin SVM

$$\underset{\mathbf{w}}{\text{argmax}} \quad \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n,m=1}^{N} \alpha_n \alpha_m t_n t_m \mathbf{x}_n^\top \mathbf{x}_m$$

$$\text{subject to} \quad \sum_{n=1}^{N} \alpha_n t_n = 0 \ \text{ and } \ 0 \le \alpha_n \le C, \text{ for all } n.$$

Optimization to find $\alpha_n$'s

$$t_{new} = \text{sign}\left( \sum_{n=1}^{N} \alpha_n t_n \mathbf{x}_n^\top \mathbf{x}_{new} + b \right)$$

predictions

*Kernelized* soft margin SVM:

$$\underset{\mathbf{w}}{\text{argmax}} \quad \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \alpha_n \alpha_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\text{subject to} \quad \sum_{n=1}^{N} \alpha_n t_n = 0 \ \text{ and } \ 0 \le \alpha_n \le C, \text{ for all } n.$$

$$t_{new} = \text{sign}\left( \sum_{n=1}^{N} \alpha_n t_n k(\mathbf{x}_n, \mathbf{x}_{new}) + b \right).$$

SISTA ⟩ 22

# Some Kernels

linear $\quad k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$

polynomial $\quad k(\mathbf{x}, \mathbf{z}) = \left(\alpha \mathbf{x}^\top \mathbf{z} + \beta\right)^\gamma$

Gaussian
(squared exponential) $\quad k(\mathbf{x}, \mathbf{z}) = \exp\left\{-\gamma \left(\mathbf{x} - \mathbf{z}\right)^\top \left(\mathbf{x} - \mathbf{z}\right)\right\}$

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \ldots$$

SISTA  23

---

# Some Kernels

linear $\quad k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$

polynomial $\quad k(\mathbf{x}, \mathbf{z}) = \left(\alpha \mathbf{x}^\top \mathbf{z} + \beta\right)$

Gaussian
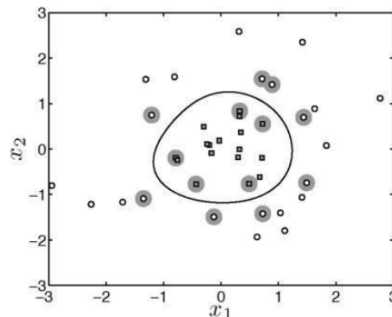(squared exponential) $\quad k(\mathbf{x}, \mathbf{z}) = \exp\left\{-\gamma \left(\mathbf{x}\right.\right.$

Kernel functions for machine learning

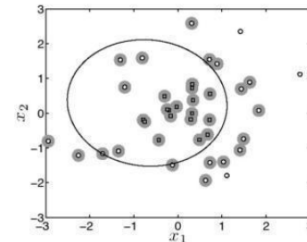http://crsouza.blogspot.com/2010/03/kernel-functions-for-machine-learning.html

1. Linear Kernel
2. Polynomial Kernel
3. Gaussian Kernel
4. Exponential Kernel
5. Laplacian Kernel
6. ANOVA Kernel
7. Hyperbolic Tangent (Sigmoid) Kernel
8. Rational Quadratic Kernel
9. Multiquadric Kernel
10. Inverse Multiquadric Kernel
11. Circular Kernel
12. Spherical Kernel
13. Wave Kernel
14. Power Kernel
15. Log Kernel
16. Spline Kernel
17. B-Spline Kernel
18. Bessel Kernel
19. Cauchy Kernel
20. Chi-Square Kernel
21. Histogram Intersection Kernel
22. Generalized Histogram Intersection Kernel
23. Generalized T-Student Kernel
24. Bayesian Kernel
25. Wavelet Kernel
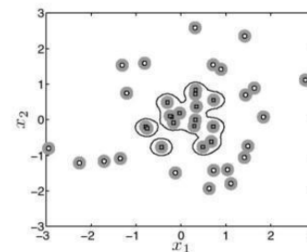
# SVM Classification with a Gaussian Kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp\left\{-\gamma(\mathbf{x}_n - \mathbf{x}_m)^\top(\mathbf{x}_n - \mathbf{x}_m)\right\}$$

(a) $\gamma = 0.01$

$\gamma = 10$, $C = 10$

$\gamma$ and $C$ are unfortunately both free parameters

(b) $\gamma = 50$

25

# SVM Odds-n-Ends

- SVM training is sensitive to **feature value ranges**.
  - This is because the view that SVM optimization has of the data is through the inner product! Greater magnitude values (positive or negative) dominate.
  - Scale your data! Common to linearly scale all feature (aka attribute) values to range [-1, 1] or [0,1].
  - Determine scaling for each feature based on your ***training data*** and **save** your scaling ranges; use the same scaling for any other data you use with the learned model.
- SVM training can be sensitive to "unbalanced" classes (one class is represented much more than another)
  - Implementations like **libsvm** allow you to weight class labels in order to "balance" the contribution of classes
- Parameter Selection:
  - Grid search: systematic combinations of parameters values and narrowing (usually doing CV at each point)

26

# Multi-class Classification

- **1-against-the-rest** :
  - Binary classifiers: class $c_i$ vs. $\{c_j, \forall j \neq i\}$.
  - Total of C classifiers created (note that that each classifier is trained on *all* data)
- **1-against-1**:
  - Binary classifiers: class $c_i$ vs. $c_j$, $j \neq i$
  - Total of $C(C\text{-}1)/2$ classifiers created (can be much faster to train if training time is super-linear in data)
- Finally, NOTE: Multi-class classification is NOT the same thing as *multi-label* **classification**: the latter involves assigning more than one label to each instance.

SISTA 27

# Kernel Nearest Neighbors

- We can Kernelize other methods

- E.g., **Nearest Neighbors**: the core of the algorithm that involves the data is a **distance metric**.

- We can turn the NN distance function into a form involving only inner products of **x**:

A common form of distance function (e.g., Euclidean distance)

$$\left(\mathbf{x}_{new} - \mathbf{x}_n\right)^{\mathsf{T}}\left(\mathbf{x}_{new} - \mathbf{x}_n\right)$$

Multiply through:  $\mathbf{x}_{new}^{\mathsf{T}}\mathbf{x}_{new} - 2\mathbf{x}_{new}^{\mathsf{T}}\mathbf{x}_n + \mathbf{x}_n^{\mathsf{T}}\mathbf{x}_n$

Kernelize:  $k(\mathbf{x}_{new}, \mathbf{x}_{new}) - 2k(\mathbf{x}_{new}, \mathbf{x}_n) + k(\mathbf{x}_n, \mathbf{x}_n)$

SISTA 28