*To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.*
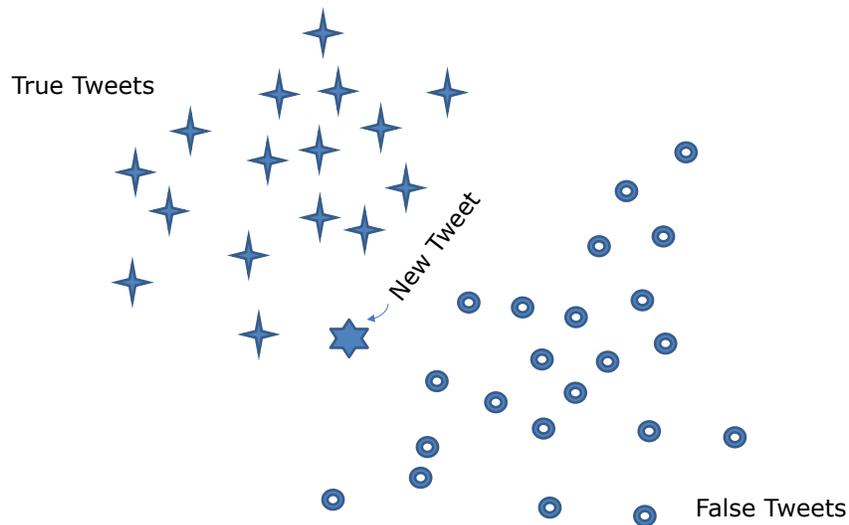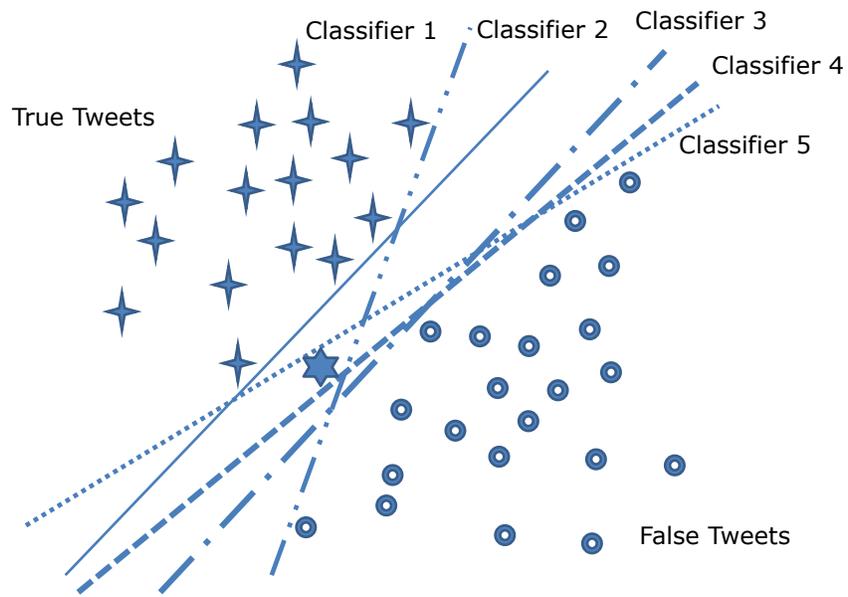
# Exploring the math in Support Vector Machines

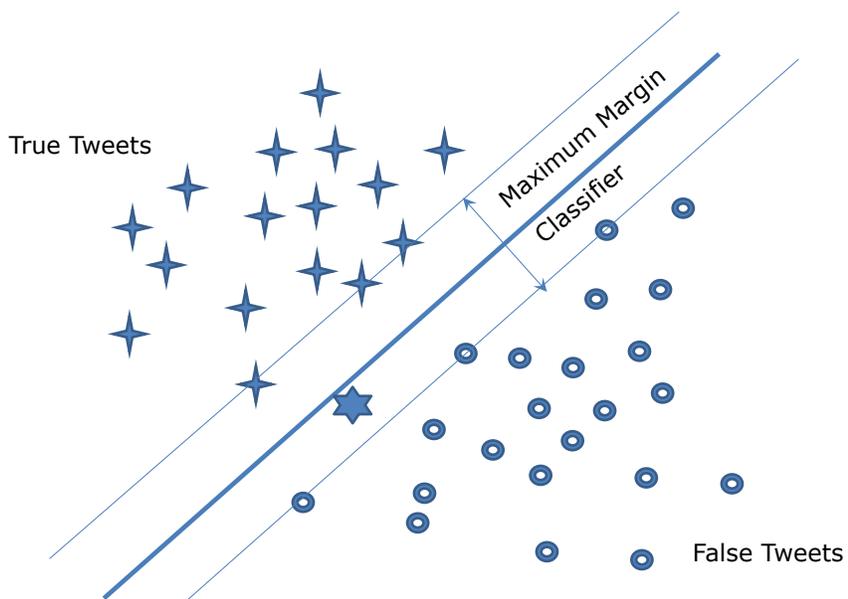VISHNU S. PENDYALA, PHD

*Video Recording:*
*https://ieeetv.ieee.org/video/exploring-the-math-in-support-vector-machines*

True Tweets

New Tweet

False Tweets

Classifier 1  Classifier 2  Classifier 3
Classifier 4
Classifier 5
True Tweets
False Tweets

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**



True Tweets
Maximum Margin
Classifier
False Tweets

Maximum Margin Classifier

Arbiter

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.



True Tweets

Support Vectors

False Tweets

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

True Tweets

Maximum Margin

Classifier

False Tweets

Vectors "supporting" the
Maximum Margin classifier

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**
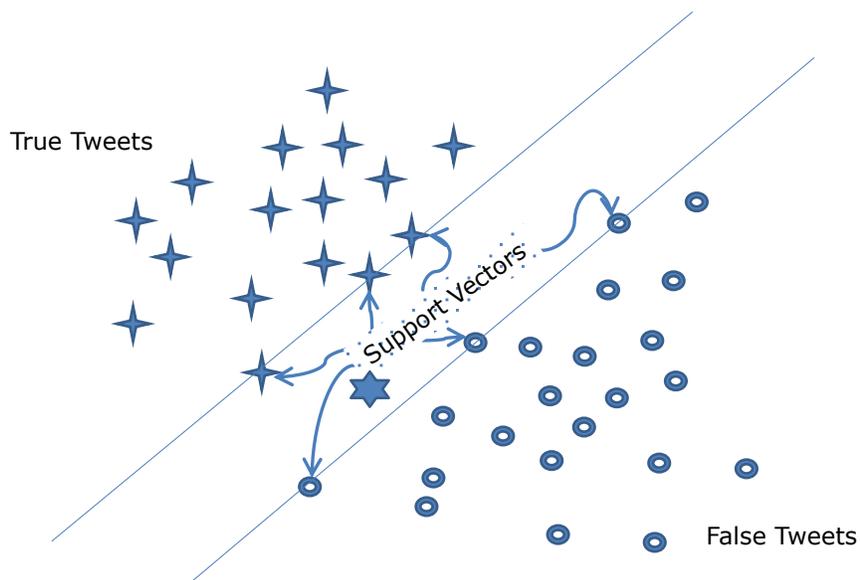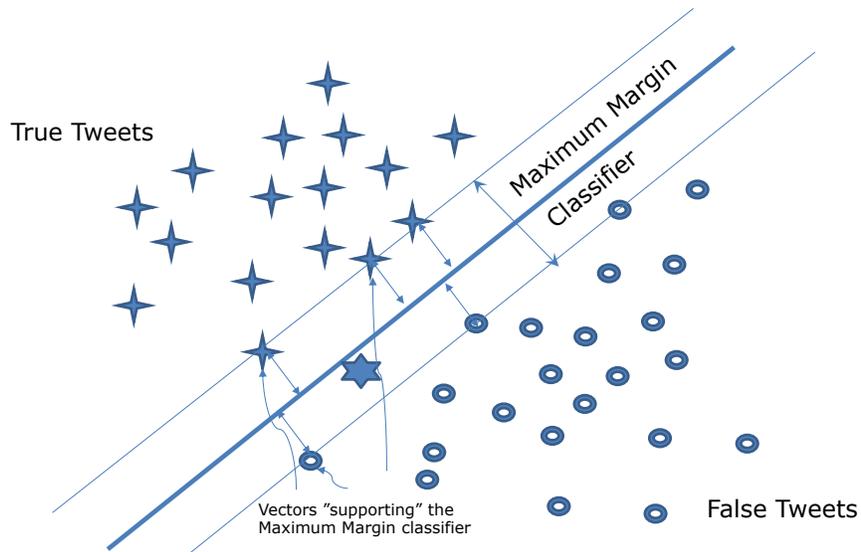
True Tweets

False Tweets

# Geometric Intuition: Convex Hulls

True Tweets

Maximum Margin

Classifier

Vectors "supporting" the
Maximum Margin classifier

False Tweets

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**
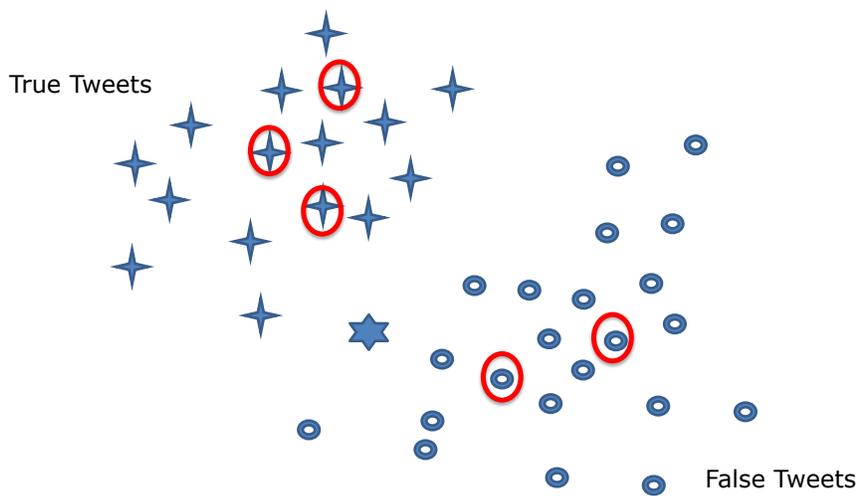
# Maximum Margin Hyperplane

$B_1$

Polytope Distance

Polytopes

$b_{11}$

$b_{12}$

## What is the length of the Maximum Margin?

Remember how to compute the distance between two parallel lines?

$$ax + by + c1 = 0$$

$$ax + by + c2 = 0$$

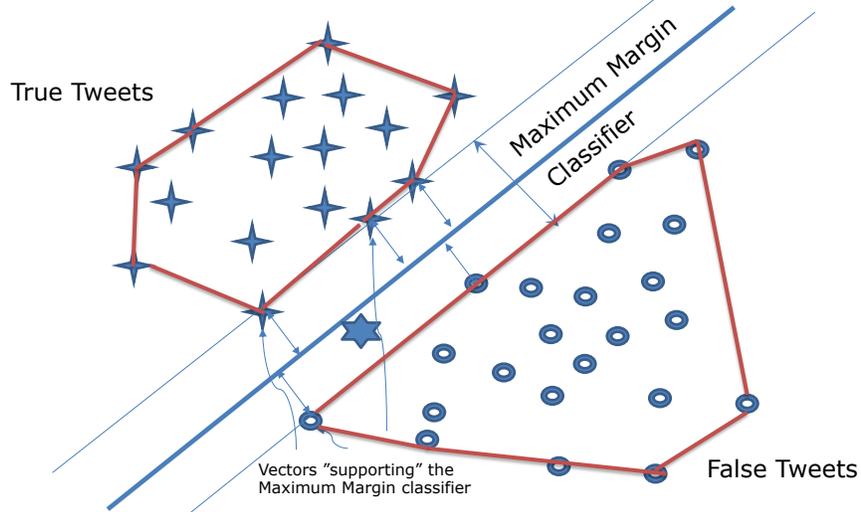$$d = \frac{|c_2 - c_1|}{\sqrt{a^2 + b^2}}$$

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

## Support Vector Classifier



$$\vec{w} \bullet \vec{x} + b = 0$$

$$\vec{w} \bullet \vec{x} + b = -1$$

$$\vec{w} \bullet \vec{x} + b = +1$$

$$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x} + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x} + b \leq -1 \end{cases}$$

$$\text{Margin} = \frac{2}{\| \vec{w} \|}$$

# Linear SVM Problem Formulation

- Objective is to maximize: $\text{Margin} = \dfrac{2}{\|\vec{w}\|}$

   a) Which is equivalent to minimizing: $L(\vec{w}) = \dfrac{\|\vec{w}\|^2}{2}$

   b) Subject to the following constraints:

$$y_i = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x}_i + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x}_i + b \leq -1 \end{cases}$$

  or

$$y_i(\mathbf{w} \bullet \mathbf{x}_i + b) \geq 1, \quad i = 1,2,\ldots,N$$

- This is a constrained optimization problem

=> Solve it using Lagrange multiplier method

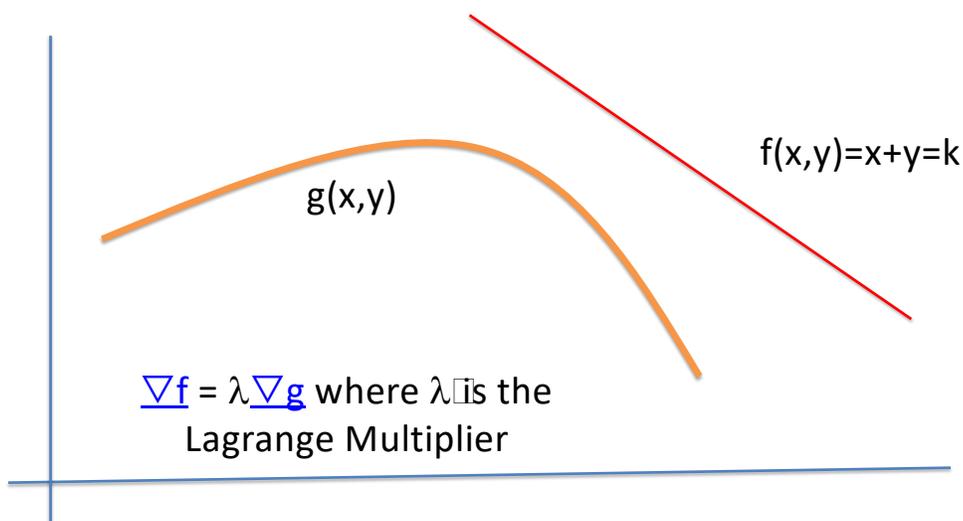# Lagrange Multipliers: Geometric Intuition

Find $\text{argmax}_{x,y} x+y$ such that $g(x,y) = 0$

f(x,y)=x+y=k

g(x,y)

$\nabla f = \lambda \nabla g$ where $\lambda$ is the Lagrange Multiplier

# Leading questions

Given two quantities f and g, how do you assign them relative importance?

Use weights: $w_0f + w_1g$

What if none of g is desirable but all of f is?

$f - w_1g$ and $w_1$ is large => $w_1g$ is a penalty term

Does this hold even when f and g are functions?

Yes! $f(x) - w_1g(x)$ can be considered an objective function that needs to be optimized if optimal x has to be determined.

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

# Constrained Optimization

1. In general, if we want to optimize f(x) under a given constraint g(x), we consider

   a) f(x) + h(g(x)) where h(g(x)) = 0 if the constraint is met and infinity if it is not met.

   b) What is a good (smoother) approximation to a step function?

$\infty$

h(g(x))

0

$h(g(x)) = \alpha g(x)$

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Lagrange Multipliers Rephrased

- Find $\text{argmax}_{x,y} f(x,y)$ such that $g(x,y)=c$

    is equivalent to finding $(x,y,\alpha)$ such that $\nabla L = 0$ where

    $L(x,y,\alpha) = f(x,y) - \alpha(g(x,y)-c)$

$\Rightarrow \delta_\alpha L(x,y,\alpha) = -g(x,y)+c=0 \ldots\ldots(1)$

$\delta_x L(x,y,\alpha) = \delta_x f(x,y) - \alpha\delta_x g(x,y)=0 \ldots(2)$

$\delta_y L(x,y,\alpha) = \delta_y f(x,y) - \alpha\delta_y g(x,y)=0 \ldots(3)$

- From (2) and (3) $\nabla f = \alpha \nabla g$

    Can also be written as

    $\delta_{x,y} f(x,y) = \alpha\delta_{x,y} g(x,y)$

Source: Abhisek Jana

# What are f and g in the case of SVM?

$$f: \quad L(\vec{w}) = \frac{\|\vec{w}\|^2}{2}$$

$$g: \quad y_i(\mathbf{w} \bullet \mathbf{x}_i + b) \geq 1, \quad i = 1,2,\ldots,N$$

How many constraints, are there in our SVM formulation?

N (look at g above: $i = 1,2,\ldots,N$)

x and w are vectors and not scalars!

What do we do now?

# The SVM Constraint is an inequality!

1. Instead of $g(x, y) = k$, we have $g(x,y) <= 0$
2. We impose Karush-Kuhn-Tucker (KKT) conditions on x, y and $\alpha$ as follows

$$L(x, y, \alpha) = f(x, y) + \alpha g(x, y) \dots..(1)$$
$$\nabla L = 0 \dots\dots\dots\dots\dots\dots\dots\dots\dots\dots(2)$$
$$g(x, y) <= 0 \dots\dots\dots\dots\dots\dots\dots\dots(3)$$
$$\alpha >= 0 \dots\dots\dots\dots\dots\dots\dots\dots.....(4)$$
$$\alpha g(x,y) = 0 \dots\dots\dots\dots\dots\dots\dots\dots.(5)$$

(5) => $\alpha = 0$ for non-support vectors; $g(x,y) = 0$ for support vectors

If there are multiple constraints, $g_i$, each will have a Lagrange multiplier $\alpha_\iota$ and the $2^{nd}$ term of L above will be a summation.

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**
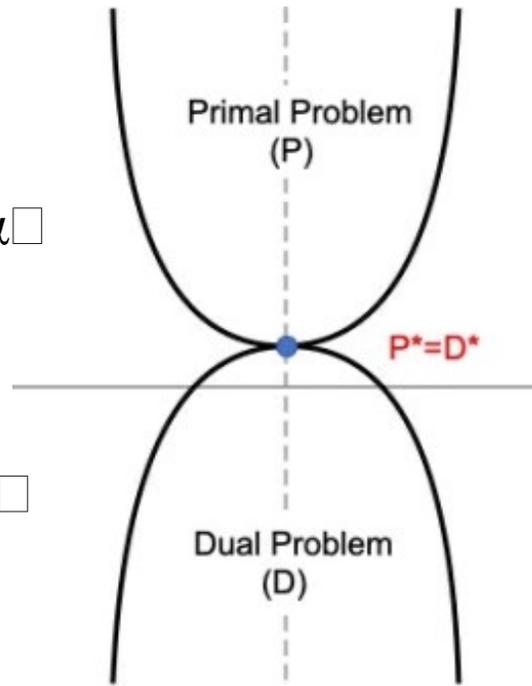
# Leading Questions

| What is the goal of x and y w.r.t. $L(x, y, \alpha)$ ? | What is the goal of $\alpha$ w.r.t. $L(x, y, \alpha)$ ? | Where did you encounter a similar scenario? | When do minmax problems stabilize? |
|---|---|---|---|
| Minimize $L(x, y, \alpha)$ | Penalty term does the opposite – maximize $L(x, y, \alpha)$ | GANs – minmax problems | Minimizing L w.r.t. x, y is same as maximizing L w.r.t. $\alpha$ => $\text{argmin}_{x,y} L(x,y,\alpha) = \text{argmax}_\alpha L(x,y,\alpha)$ |

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Strong Duality

$$\min{}_{w,b} \; L(w, b, \alpha)$$

Primal Problem
(P)

$P*=D*$

$$\max_{\alpha} \; L(w, b, \alpha)$$

Dual Problem
(D)

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

# Leading Questions

What is the abstraction for something that can take a finite number of values?

Vector: [x1, x2, x3,…, xN]

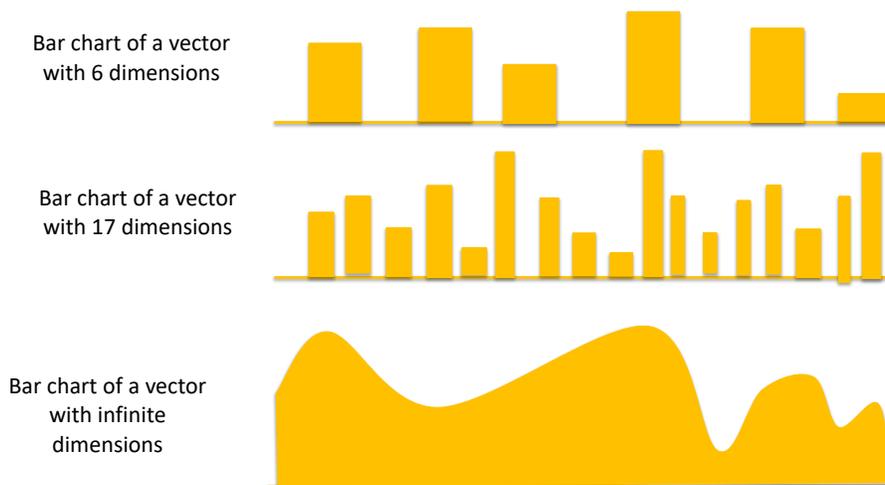What is the abstraction for something that can take infinite number of values?

A variable: x

What is it that can take infinite number of values that obey certain properties?

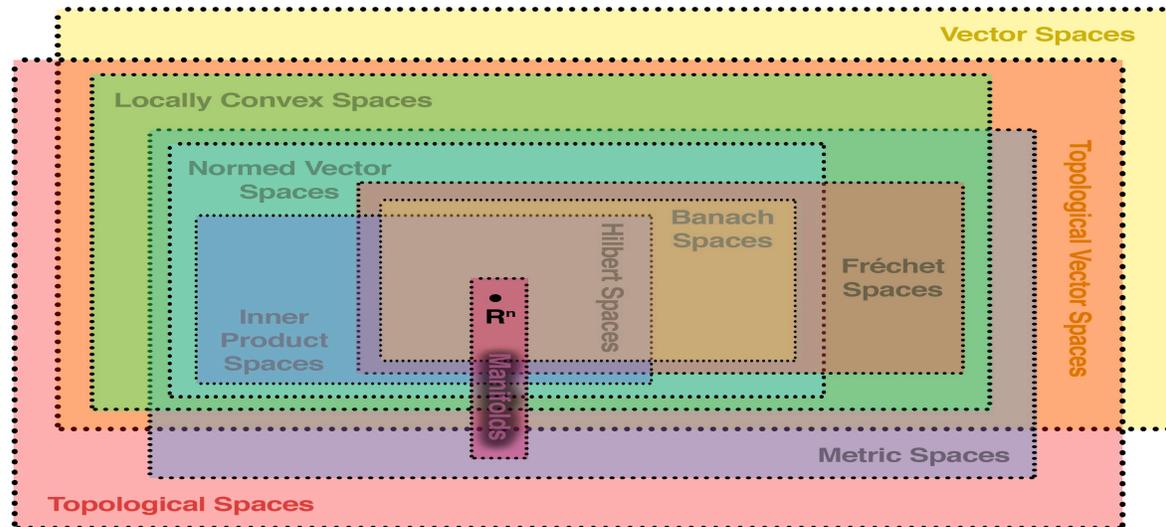Function y = f(x)

Functions are infinite dimensional vectors!

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Remember discrete vs continuous random variables and PMF / PDF?

Bar chart of a vector with 6 dimensions

Bar chart of a vector with 17 dimensions

Bar chart of a vector with infinite dimensions

| Abstraction | Vector Space | Function Space |
|---|---|---|
| **Summation** | $\Sigma$ | $\int$ |
| **Norm** | $\parallel x \parallel_p = \left(\Sigma |x|^p\right)^{1/p}$ | $\parallel f \parallel_p = \left(\int |f(x)|^p dx\right)^{1/p}$ |
| **Inner Product** | $<u, v> = \Sigma u_i . v_i$ | $<f, g> = \int f(t).g(t)dt$ |
| **Orthogonality** | $<u, v> = 0$ | $<f, g> = 0$ |
| **Bases** | Orthonormal vectors | Orthonormal functions |

# The Dual is in Hilbert Space

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

## Solving for w, b, and $\alpha$: Part-1: Primal (w and b)

$L(x,y,\alpha) = f(x,y) + \alpha(g(x,y) - c)$

$L(w, b, \alpha) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{n} \alpha_i(1 - y_i(w^T x_i + b))$

$$\nabla L = 0 \Rightarrow \delta_w L = w - \sum_{i=1}^{n} \alpha_i y_i x_i = 0 \Rightarrow w = \sum_{i=1}^{n} \alpha_i y_i x_i$$

**Classifier**: $h(x) = \text{sgn}(w \cdot x + b) = \text{sgn}(\sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x_i} \boldsymbol{x} + b)$

For support vectors, $w.x_j + b = y_j$

$\Rightarrow b = y_j - \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x_i} \boldsymbol{x_j}$

$$\nabla L = 0 \Rightarrow \delta_b L = \sum_{i=1}^{n} \alpha_i y_i = 0$$

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Solving for w, b, and $\alpha$: Part-2: Dual ($\alpha$)

$L(w, b, \alpha) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{n}\alpha_i(1 - y_i(w^T x_i + b))$

$$= \frac{1}{2}w^T w - \sum_{i=1}^{n}\alpha_i y_i w^T x_i - \sum_{i=1}^{n}\alpha_i y_i b + \sum_{i=1}^{n}\alpha_i$$

*Substitute* $w = \sum_{i=1}^{n}\alpha_i y_i x_i$ *and* $\sum_{i=1}^{n}\alpha_i y_i = 0$ *from previous slide*

$$L = \frac{1}{2}w^T w - w^T w - b(0) + \sum_{i=1}^{n}\alpha_i = -\frac{1}{2}w^T w + \sum_{i=1}^{n}\alpha_i$$

$= \sum_{i=1}^{n}\alpha_i - \frac{1}{2}(\sum_{i=1}^{n}\alpha_i y_i x_i)(\sum_{j=1}^{n}\alpha_j y_j x_j)$

$= \sum_{i=1}^{n}\alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i \alpha_j y_i y_j (x_i^T x_j)$

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.
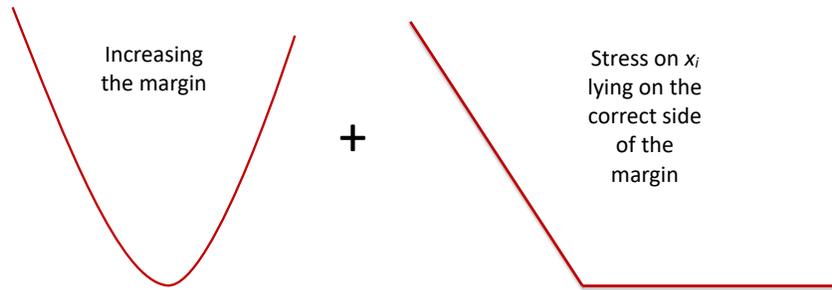
## Next Steps

- We use Quadratic Programming software or algorithms like SMO (Sequential Minimal Optimization) to solve for alphas
- At optimal solution, the data items corresponding to non-zero alphas are the support vectors
- Support vectors and the corresponding alphas can be saved as a model
- The feature vectors in the training set appear only inside dot products => ready for the kernel trick for non-linear data!

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

# Notes

1. We **converted** a problem in terms of the weight vector $w$ and bias $b$ into a problem in terms of the Lagrange multipliers $\alpha_i$

2. Feature vectors $x_i$ and labels $y_i$ are **already known** from the training dataset.

3. Since the variables in our problem changed resulting in a dual, instead of minimizing L w.r.t. w and b , **we must now maximize L w.r.t** $\alpha_i$

4. The original problem is actually equivalent to $\min_{w,b} \max_{\alpha}$ L(w, b, $\alpha$)

= $\min_{w,b} \max_{\alpha}$ (f(w, b) + $\alpha$g(w, b))

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.
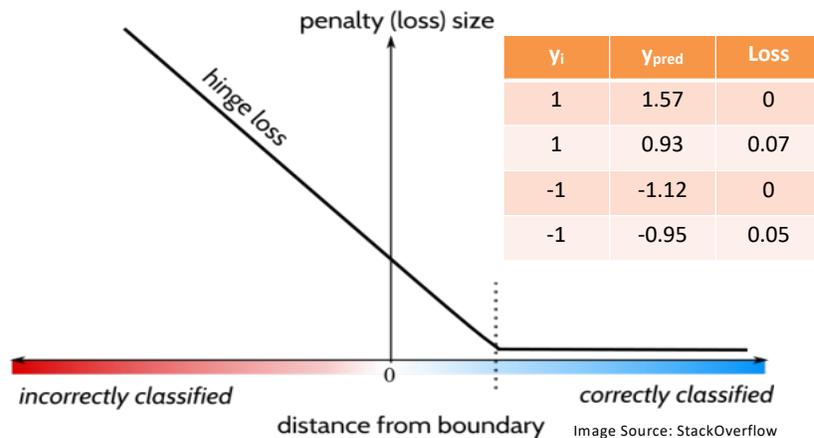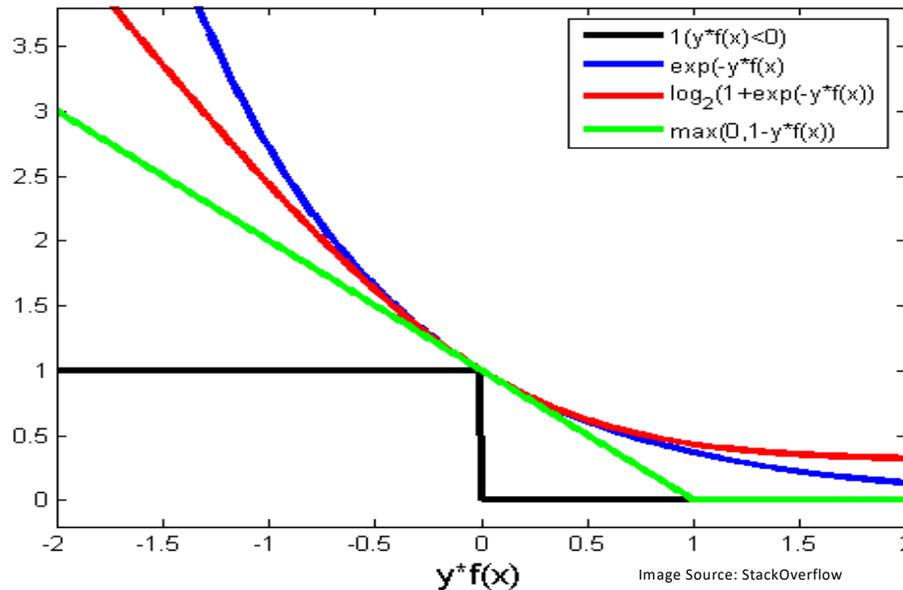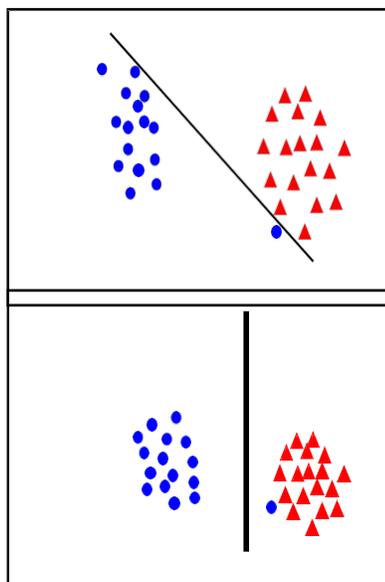
# Observation

- When we use Lagrange "trick", QP, and duality, we notice that all the equations have only **dot products** of the data instances

- Nowhere will we need the individual values of the features

- Classification is anyway all about **pairwise similarity** of the data points, which the dot product indicates

- This finding comes in **handy** when we consider **non-linear** classification

# SVM Objective Function is Convex

$$\arg\min_w \frac{1}{2}\|w\|^2 + \sum_{i=1}^{n} \alpha_i \max[0, 1 - y_i w^T x_i]$$

Regularization                    (Hinge) Loss Function

Increasing the margin

**+**

Stress on $x_i$ lying on the correct side of the margin

Sum of Convex Functions is Convex

# Hinge Loss

$$\arg\min_w \frac{1}{2}\|w\|^2 + \sum_{i=1}^{n} \alpha_i \max[0, 1 - y_i w^T x_i]$$

Regularization                    (Hinge) Loss Function

penalty (loss) size

hinge loss

| $y_i$ | $y_{pred}$ | Loss |
|-------|------------|------|
| 1     | 1.57       | 0    |
| 1     | 0.93       | 0.07 |
| -1    | -1.12      | 0    |
| -1    | -0.95      | 0.05 |

0

*incorrectly classified*                    *correctly classified*

**distance from boundary**    Image Source: StackOverflow

# Various Loss Functions



Image Source: StackOverflow

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

## Linear separability: What is the best w?



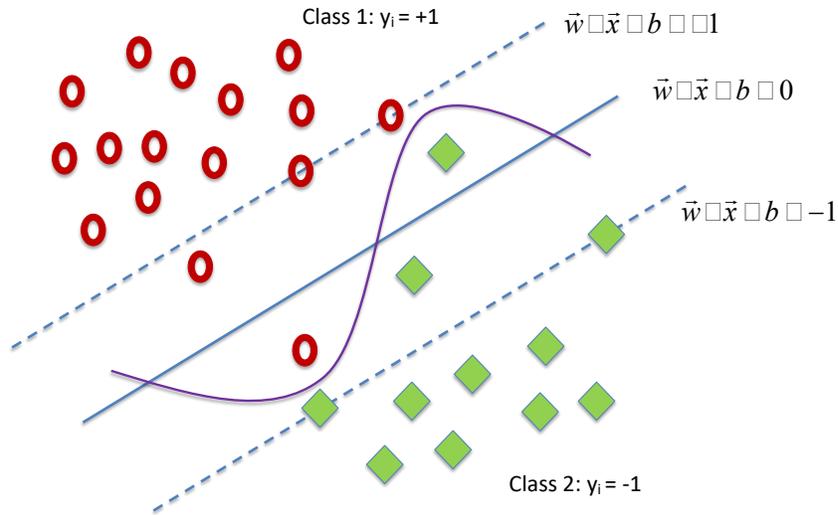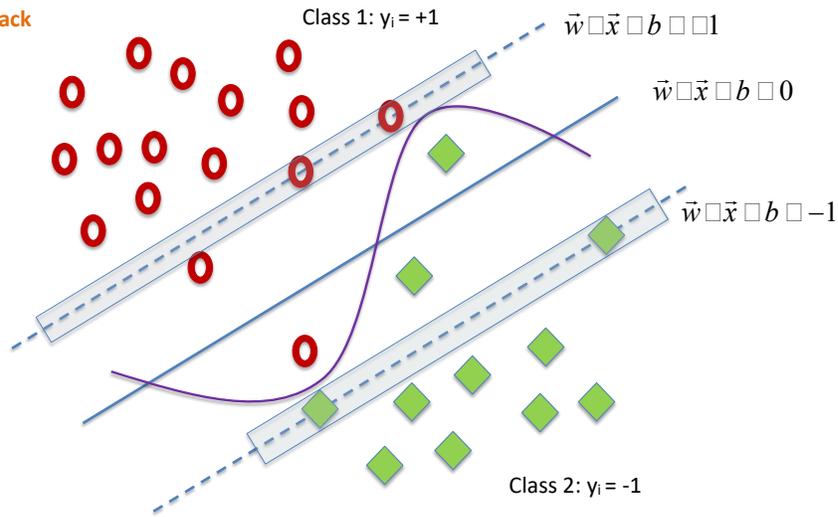• the points can be linearly separated but there is a very **narrow margin**

• but possibly the large margin solution is better, even though **one constraint is violated**

In general there is a trade off between the margin and the number of mistakes on the training data

Source: A. Zisserman

# Slightly non-linear: Can we cut some slack?

Class 1: $y_i = +1$

$\vec{w} \bullet \vec{x} + b = +1$

$\vec{w} \bullet \vec{x} + b = 0$

$\vec{w} \bullet \vec{x} + b = -1$

Class 2: $y_i = -1$

We keep track of the "bad apples" using "slack" variables

11/12/23

# Category #1: Correctly Classified Points

$y_i(W.X+b) > 1$

All Good – no need for slack

Class 1: $y_i = +1$

$\vec{w} \bullet \vec{x} + b = +1$

$\vec{w} \bullet \vec{x} + b = 0$

$\vec{w} \bullet \vec{x} + b = -1$

Class 2: $y_i = -1$

11/12/23

# Category #2: Support Vectors

$y_i(W.X+b) = 1$

**All Good – no need for slack**

Class 1: $y_i = +1$

$\vec{w} \bullet \vec{x} + b = +1$

$\vec{w} \bullet \vec{x} + b = 0$

$\vec{w} \bullet \vec{x} + b = -1$

Class 2: $y_i = -1$

# Category #3: Wrong Side of Margin

**But classified correctly**

$y_i(W.X+b) < 1$

$0 < \zeta_i \text{ (slack)} < 1$

Class 1: $y_i = +1$

$\vec{w} \bullet \vec{x} + b = +1$

$\vec{w} \bullet \vec{x} + b = 0$

$\vec{w} \bullet \vec{x} + b = -1$

Class 2: $y_i = -1$

# Category #4: Incorrectly Classified
### …and wrong side of margin

$y_i(W.X+b) < 1$

$\zeta_i > 1$

Class 1: $y_i = +1$

$\vec{w} \bullet \vec{x} + b = +1$

$\vec{w} \bullet \vec{x} + b = 0$

$\vec{w} \bullet \vec{x} + b = -1$

Class 2: $y_i = -1$

11/12/23

# Soft Margin SVM: Objective Function

We add a slack term and a hyperparameter, C:

$$\text{Minimize } \frac{\|W\|^2}{2} + C \sum \zeta_i$$

and the constraints change as well

$$y_i(W.X + b) >= 1 - \zeta_i$$

Repeat the math with this additional set of variables, $\zeta_i$ – Gradients now include w.r.t $\zeta_i$

- Large C => close to hard margin, less points misclassified, smaller margin, overfitting
- C-->infinity => hard margin

# Hyperparameter C and the trade-off

$$\arg\min_w \frac{1}{2}\|w\|^2 \quad + \quad C\sum_{i=1}^{n}\xi_i$$

| Increasing the margin | Vs | Stress on $x_i$ lying on the correct side of the margin |
|---|---|---|

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

# Non-linear SVM

Data cannot be separated by a single hyperplane

# What do we do when our 1-bed apartment gets noisy and we can't work?

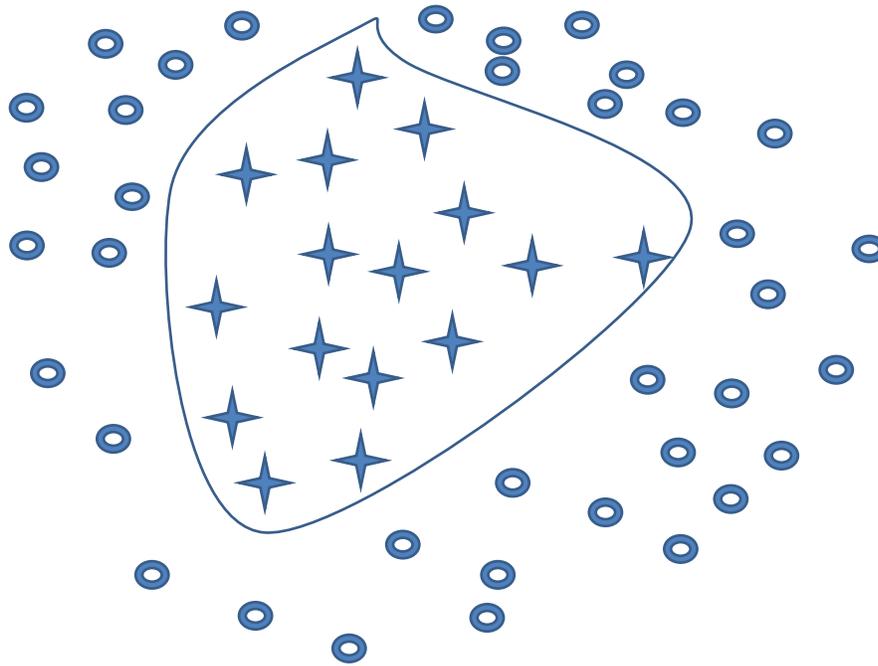We go to the library which is more spacious and quieter!

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

# What if there was trick that could transport you from the 3D world to a space with infinite dimensions?

Or better still: Let you work in infinite dimensions without even visiting the space?

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

## Nonlinear SVMs

- Linearly separable dataset in 1D:



- Non-separable dataset in 1D:



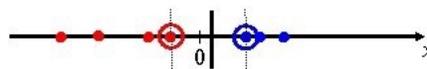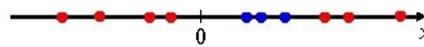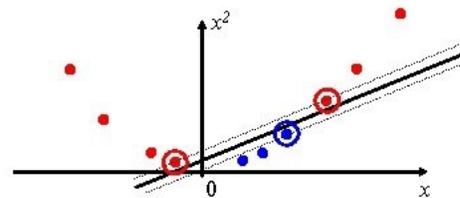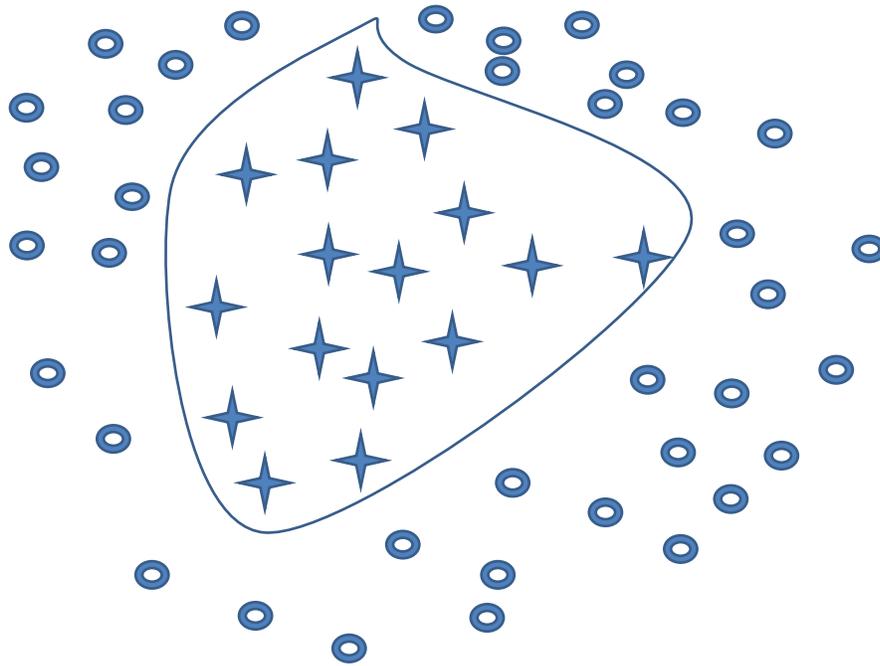- We can map the data to a *higher-dimensional space*:

Slide credit: Andrew Moore

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**
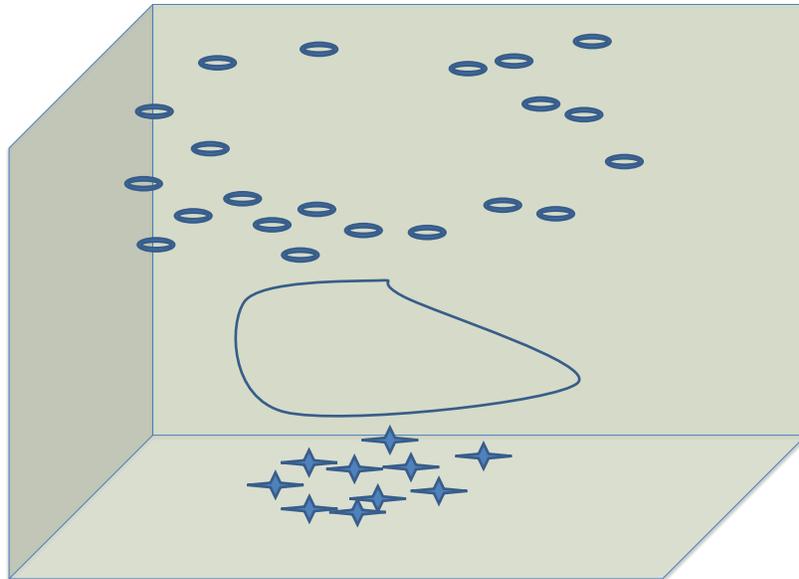
# WHEN MAPPED TO 3-D…

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.
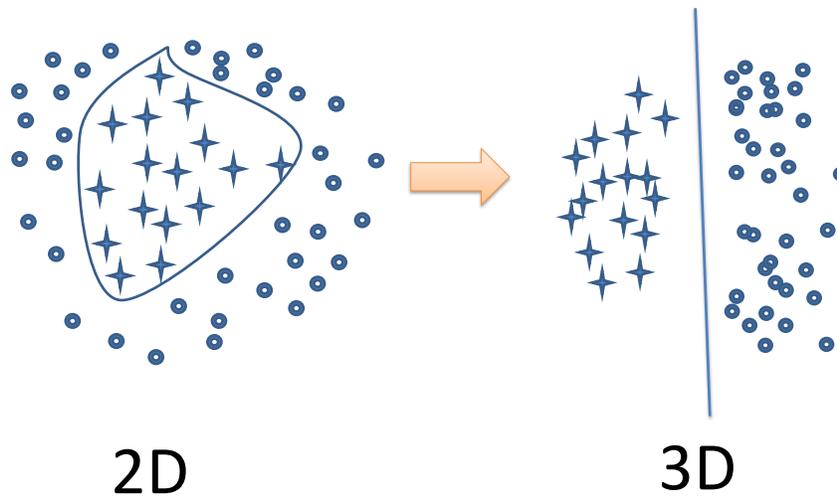
To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

# Map from 2D to 3D



2D                    3D

# Cover's Theorem

"Given a set of training data that is not linearly separable, one can, with high probability transform it into a training set that is linearly separable by projecting it into a higher-dimensional space via some non-linear transformation." [Proof](#)

**Thomas M. Cover**
Photo credit: IEEE Information Theory Society

## Cover's theorem illustrated

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \xrightarrow{\phi} \quad \begin{pmatrix} z_1 = x_1^2 \\ z_2 = \sqrt{2}x_1 x_2 \\ z_3 = x_2^2 \end{pmatrix}$$

$\phi$

$\phi$ maps from 2D to 3D, making the problem linearly separable

# What is $\phi$?

Mapping Function:

$$\phi: \quad \begin{aligned} \Re^2 &\longrightarrow \Re^3 \\ (x_1, x_2) &\longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{aligned}$$

Equation of the Classifying Hyperplane in 3D:

$$\omega^T z + b = 0$$

Substituting for z, we get the equation of an ellipse in 2D:

$$\omega_1 x_1^2 + \omega_2 \sqrt{2} x_1 x_2 + \omega_3 x_2^2 = 0$$

$\phi$ Maps x $\in X$ $to$ $R^D$ where D can be potentially infinite
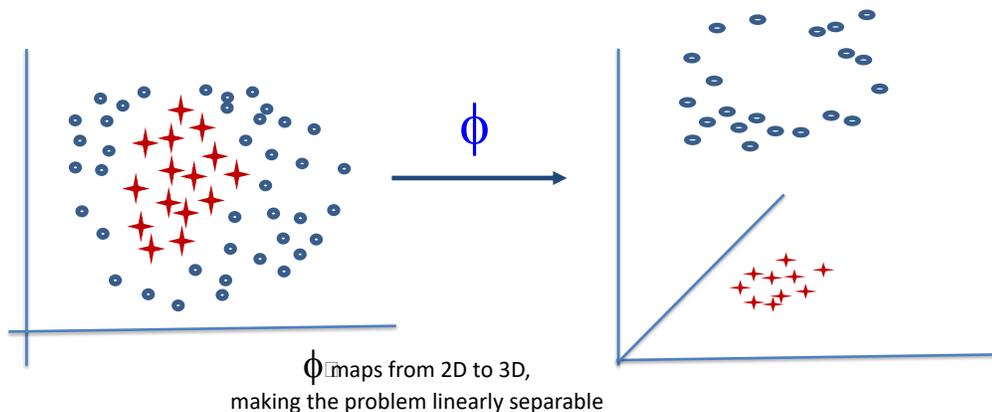
Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

"The ordinary operations of algebra suffice to resolve problems in the theory of curves.

As long as algebra and geometry have been separated, their progress have been slow and their uses limited; but when these two sciences have been united, they have lent each mutual forces, and have marched together towards perfection." - Joseph-Louis Lagrange

**Born:** January 25, 1736
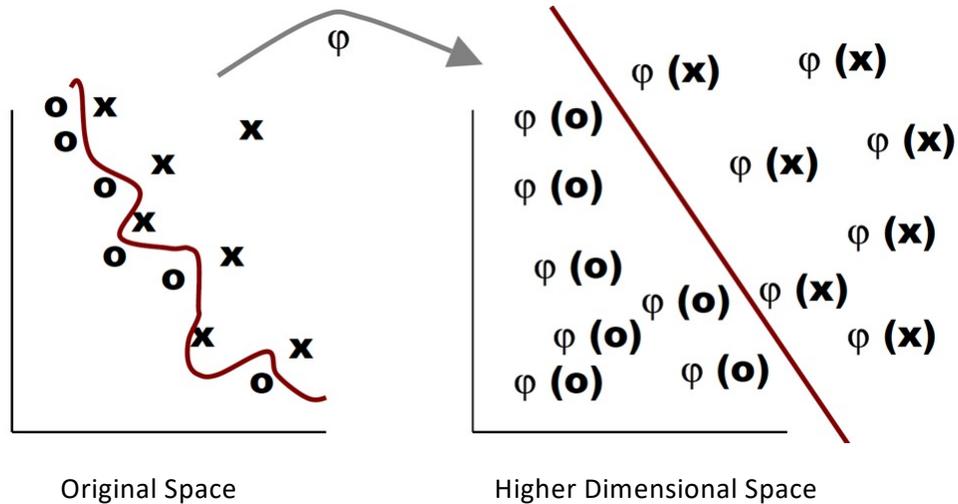**Died:** April 10, 1813
Photo credit: Wikipedia

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

## A different perspective of Cover's theorem

1. Machine Learning is essentially is fitting a curve $y = f(x)$ to the data.
2. We fit a straight line if the dependencies among the data are linear, $y = mx + c$
3. How can we fit a straight line if the relationships in the data are non-linear, say modeled by $y = f(x) = mx^2 + c$?
4. Export the data to a different space where $x' = mx^2 + c$ $then$ $y' = x'$ fits the data!
5. SVM finds $y = f(x)$ by transforming the data to a different space.

# Transformation to Higher Dimension



Original Space                    Higher Dimensional Space

Image Source: R. Berwick

# Problems with the transformation

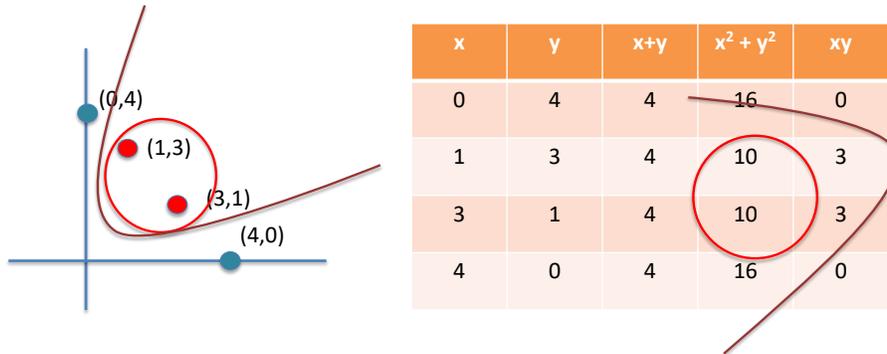| | |
|---|---|
| ✔ | If there are 1M training samples, need to evaluate the function $\phi$ at 1M points |
| 📐 | If the transformed space has 1M dimensions, need to evaluate dot products across all the 1M dimensions for all pairs of samples |
| 🧠 | => Computationally infeasible |

# Not all transformations apply!

| x | y | x+y | $x^2 + y^2$ | xy |
|---|---|-----|-------------|----|
| 0 | 4 | 4 | 16 | 0 |
| 1 | 3 | 4 | 10 | 3 |
| 3 | 1 | 4 | 10 | 3 |
| 4 | 0 | 4 | 16 | 0 |

(0,4)
(1,3)
(3,1)
(4,0)

## Which transformations apply and how do we solve all these problems?

11/12/23

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

Lesson: Capture only the core (kernel) and leave the rest!

How long does it take to train a model in self-driving car to recognize objects?

Only dot products are needed in the dual formulation

11/12/23

## A "Gram" Matrix of inner products in the transformed space is all we need

| $< \phi(x_1), \phi(x_1) >$ | $< \phi(x_1), \phi(x_2) >$ | ... | $< \phi(x_1), \phi(x_n) >$ |
|---|---|---|---|
| $< \phi(x_2), \phi(x_1) >$ | $< \phi(x_2), \phi(x_2) >$ | ... | $< \phi(x_2), \phi(x_n) >$ |
| ... | ... | ... | ... |
| $< \phi(x_n), \phi(x_1) >$ | $< \phi(x_n), \phi(x_2) >$ | ... | $< \phi(x_n), \phi(x_n) >$ |

## Simplifying the Gram matrix computation

- How can we avoid the computing cost of the (i) transformations and (ii) inner products?
- What if we found a function

$$K(x_m, x_n) = < \phi(x_m), \phi(x_n) > \forall x_m, x_n$$

- We can compute K in the original feature space without computing $\phi's$ or their inner products
- $K$ should satisfy the properties of a Gram matrix
- Classifier is given by

$$f(x) = \text{sgn}(\sum_{i=1}^{n} \alpha_i y_i \phi(x_i)^T \phi(x)) + b = \text{sgn}(\sum_{i=1}^{n} \alpha_i y_i K(x_i, x)) + b$$

# Kernel Gram matrix as a function

| $\mathbf{K}$ | 1 | 2 | $\cdots$ | $\ell$ |
|---|---|---|---|---|
| 1 | $\kappa\left(\mathbf{x}_1, \mathbf{x}_1\right)$ | $\kappa\left(\mathbf{x}_1, \mathbf{x}_2\right)$ | $\cdots$ | $\kappa\left(\mathbf{x}_1, \mathbf{x}_\ell\right)$ |
| 2 | $\kappa\left(\mathbf{x}_2, \mathbf{x}_1\right)$ | $\kappa\left(\mathbf{x}_2, \mathbf{x}_2\right)$ | $\cdots$ | $\kappa\left(\mathbf{x}_2, \mathbf{x}_\ell\right)$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $\ell$ | $\kappa\left(\mathbf{x}_\ell, \mathbf{x}_1\right)$ | $\kappa\left(\mathbf{x}_\ell, \mathbf{x}_2\right)$ | $\cdots$ | $\kappa\left(\mathbf{x}_\ell, \mathbf{x}_\ell\right)$ |

As $\ell \to \infty$, the above infinite dimensional matrix can be represented by a function under certain circumstances

First, we need an (i) an inner product space that is (ii) complete to which the function maps

# Kernel Functions

- Represent inner products in the transformed space => need to satisfy the properties of inner product, such as symmetry $< x_1, x_2 > = < x_2, x_1 >$ and positive semi-definiteness $< x, x > \geq 0$

- Implies that the Gram matrix has to be symmetric and positive semi-definite (=> Eigenvalues of the matrix are non-negative)

# Mercer Conditions

For a kernel function to truly map to an inner product, the kernel gram matrix must also be

1. Symmetric => $K(x, y) = K(y, x)$ and

2. The diagonal elements represent the squared norms, therefore must be positive

3. Positive semidefinite =>

$\Sigma_i \Sigma_j \alpha_i \alpha_j K(x_i, x_j) >= 0$ for all $\alpha_i$ and $\alpha_j$ in R

(or) $v^T K v \geq 0 \ \forall v$

(or) its eigenvalues are nonnegative. Why?

# Why positive semi-definite?

- Suppose $\phi(x)$ = x, i.e., no transformation
- $K = X^T X$ where X is the training dataset features
- Consider any vector $v$ with elements $\alpha_i$
$$v^T K v = v^T X^T X v = (Xv)^T (Xv) = u^T u \geq 0$$
where $u$ is a new vector; the above is same as
$\Sigma_i \Sigma_j \alpha_i \alpha_j K(x_i, x_j) >= 0$
- Aside: The logic applies to covariance matrix as well

# Simplified Mercer's Theorem

"Let $K: \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$ be given.

Then for K to be a valid (Mercer) kernel, it is necessary and sufficient that for any $\{x^{(1)}, \dots, x^{(m)}\}, (m < \infty),$

the corresponding kernel matrix is symmetric positive semi-definite."

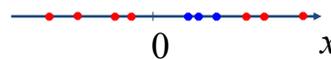*For Mercer kernels, there exists a transformation function $\phi$ such that*

*$K(x,x') = \phi(x)^T \cdot \phi(x')$*

## The Kernel Trick

- Non-separable data in 1D:



- Apply mapping $\varphi(x) = (x, x^2)$:

Think of kernel as a function that directly gives you a similarity metric in the transformed space



- Dot Product of the new, transformed vectors
  - $\varphi(x)^T \varphi(x') = K(x, x') = xx' + x^2 x'^2$

## Kernel example 1: Polynomial

- Polynomial kernel with degree $d$ and constant $c$:
$$K(x, x') = (x^T x' + c)^d$$
- What this looks like for $d = 2$:
$$x = (u, v), \qquad x' = (u', v')$$
$$K(x, x') = (uu' + vv' + c)^2$$
$$= u^2 u'^2 + v^2 v'^2 + 2uu'vv' + 2cuu' + 2cvv' + c^2$$
$$= \phi(x)^\top. \; \phi(x), where$$
$$\phi(x) = (u^2, v^2, \sqrt{2}uv, \sqrt{2cu}, \sqrt{2cv}, c)$$
  We mapped 2D to 6D space!
- Thus, the explicit feature transformation consists of all polynomial combinations of individual dimensions of degree up to $d$

  11/12/23

# What if the data is extremely nonlinear?



C = 7.9
Using Rbf kernel with sigma = 0.50
Number of support vectors: 20 / 20
Converged in 41 iterations.

Source: Andrej Karpathy

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

RBF Kernel
Intuition

11/12/23

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.
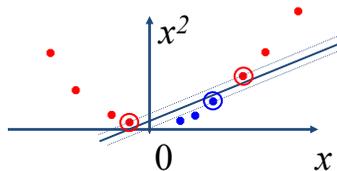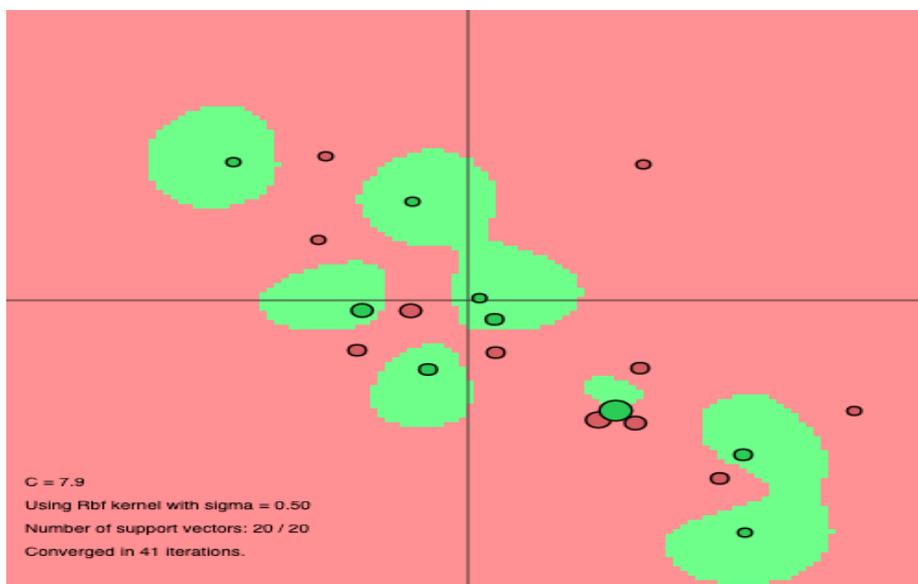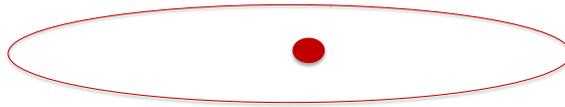
To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

The value of the RBF
Kernel is maximum at
the center where the
distance = 0
$e^0 = 1$

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

11/12/23

**RBF Kernel Space**

1. Gaussian radial basis function (RBF) kernel: $K(X_i, X_j) =$
$$e^{-\|X_i - X_j\|^2 / 2\sigma^2}$$

1. Suppose there are 5 original 2-dimensional points:
   a) $x_1(0, 0)$, $x_2(4, 4)$, $x_3(-4, 4)$, $x_4(-4, -4)$, $x_5(4, -4)$
2. If we set $\sigma$ to 4, we will have the following points in the kernel space
   a) E.g., $\|x_1 - x_2\|^2 = (0 - 4)^2 + (0 - 4)^2 = 32$, thus, $K(x_1, x_2) = e^{-\frac{32}{2 \cdot 4^2}} = e^{-1}$
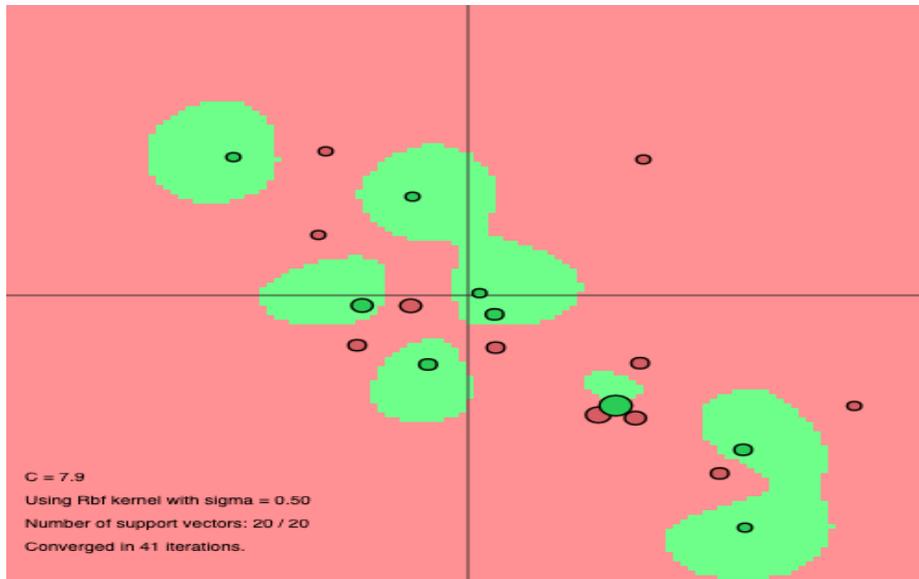
| Original Space | | |
|---|---|---|
| | $x$ | $y$ |
| $x_1$ | 0 | 0 |
| $x_2$ | 4 | 4 |
| $x_3$ | −4 | 4 |
| $x_4$ | −4 | −4 |
| $x_5$ | 4 | −4 |

RBF Kernel Space ($\sigma = 4$)

| $K(x_i, x_1)$ | $K(x_i, x_2)$ | $K(x_i, x_3)$ | $K(x_i, x_4)$ | $K(x_i, x_5)$ |
|---|---|---|---|---|
| 0 | $e^{-\frac{4^2 + 4^2}{2 \cdot 4^2}} = e^{-1}$ | $e^{-1}$ | $e^{-1}$ | $e^{-1}$ |
| $e^{-1}$ | 0 | $e^{-2}$ | $e^{-4}$ | $e^{-2}$ |
| $e^{-1}$ | $e^{-2}$ | 0 | $e^{-2}$ | $e^{-4}$ |
| $e^{-1}$ | $e^{-4}$ | $e^{-2}$ | 0 | $e^{-2}$ |
| $e^{-1}$ | $e^{-2}$ | $e^{-4}$ | $e^{-2}$ | 0 |

Source: Jiawei Han

# What if the data is extremely nonlinear?



C = 7.9
Using Rbf kernel with sigma = 0.50
Number of support vectors: 20 / 20
Converged in 41 iterations.

Source: Andrej Karpathy

# Important Kernel Functions

| Name | Kernel function | dim $(\mathcal{K})$ |
|---|---|---|
| $p$th degree polynomial | $k\,(\vec{u}, \vec{v}) = (\langle \vec{u}, \vec{v} \rangle_{\mathcal{X}})^p$ <br> $p \in \mathbb{N}^+$ | $\binom{N+p-1}{p}$ |
| complete polynomial | $k\,(\vec{u}, \vec{v}) = (\langle \vec{u}, \vec{v} \rangle_{\mathcal{X}} + c)^p$ <br> $c \in \mathbb{R}^+, \ p \in \mathbb{N}^+$ | $\binom{N+p}{p}$ |
| RBF kernel | $k\,(\vec{u}, \vec{v}) = \exp\left(-\dfrac{\|\vec{u} - \vec{v}\|_{\mathcal{X}}^2}{2\sigma^2}\right)$ <br> $\sigma \in \mathbb{R}^+$ | $\infty$ |

Source: R. Herbrich

# How Many Dimensions?

1. Polynomial Kernel when d is 2 and # of features is 2 (previous slide) results in a new feature space of 6 dimensions

2. In general, a polynomial kernel will result in a new feature space with $^{n+p}C_p$ dimensions

3. How about for RBF? Think Taylor's series:

$$\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!}(x - a)^n \qquad \Rightarrow \text{ infinite dimensions}$$

$e^x = 1 + x + x^2 / 2! + x^3 / 3! + x^4 / 4! + x^5 / 5! + \ldots$ at a = 0

RBF Kernel as a sum of Polynomial Kernels

Follows from Taylor Series:
$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots, \quad -\infty < x < \infty$$

=> Infinite dimensions

11/12/23

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**

# The Kernel Trick

1. Replace inner products in original feature space, <x,x'> with a call to the kernel function, K(x,x')

2. K(x,x') can be written as a dot product of transformation function $\phi$ as $\phi(x)^T . \phi(x')$

3. In effect, we have transformed

$$<x,x'> \longrightarrow \ <\phi(x), \phi(x')>$$

=> Transforming the data in original feature space, X to data in a transformed feature space in higher dimensions

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

“SVMs are a rare example of a methodology where geometric intuition, elegant mathematics, theoretical guarantees, and practical algorithms meet” – Bennet and Campbell
Bennett, Kristin P., and Colin Campbell. "Support vector machines: hype or hallelujah?." *ACM SIGKDD Explorations Newsletter* 2.2 (2000): 1-13.

# Relating SVM to other ideas in Machine Learning

# Kernel SVM and K-NN

$$f(x) = \text{sgn}(\sum_{i=1}^{n} \alpha_i y_i \phi(x_i)^T \phi(x)) = \text{sgn}(\sum_{i=1}^{n} \alpha_i y_i K(x_i, x))$$

- If f(x) > 0, y = +1 otherwise, y = -1

- K is an inner product, which is a distance metric

- $\alpha_i$ = 0 for non-support vectors, by KKT

- $\alpha_i$ can be viewed as a weight or relative importance of each support vector

- This is instance-based learning, like K-NN
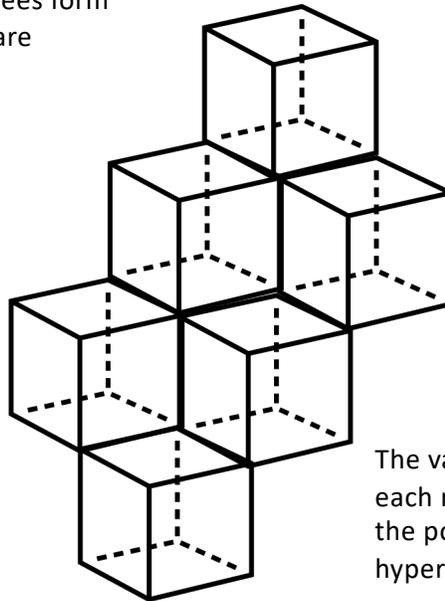
- Here the instances are support vectors

Some or all of the slides in this presentation may have been influenced by or adopted from various sources for the sole purpose of teaching students and enhancing their learning experience.

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.



SVM and Regression Trees

Regression Trees form regions that are <mark>hypercubes</mark>

Hypercubes are formed based on <mark>decision splits</mark>



The value associated with each region is $Y_{average}$ of the points in that hypercube

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.

# Single Feature - Binning

Goal is to find the splits that minimize MSE



Preferred Customer Ranking

income

# Regression Trees

- In regression where y is continuous, a test item's expected target variable is predicted as

$$f(x) = \sum_{m=1}^{M} \left( \frac{1}{k} \left( \sum_{i=1}^{k} y_i \right) \mathbb{I}(X \in Rm) \right)$$
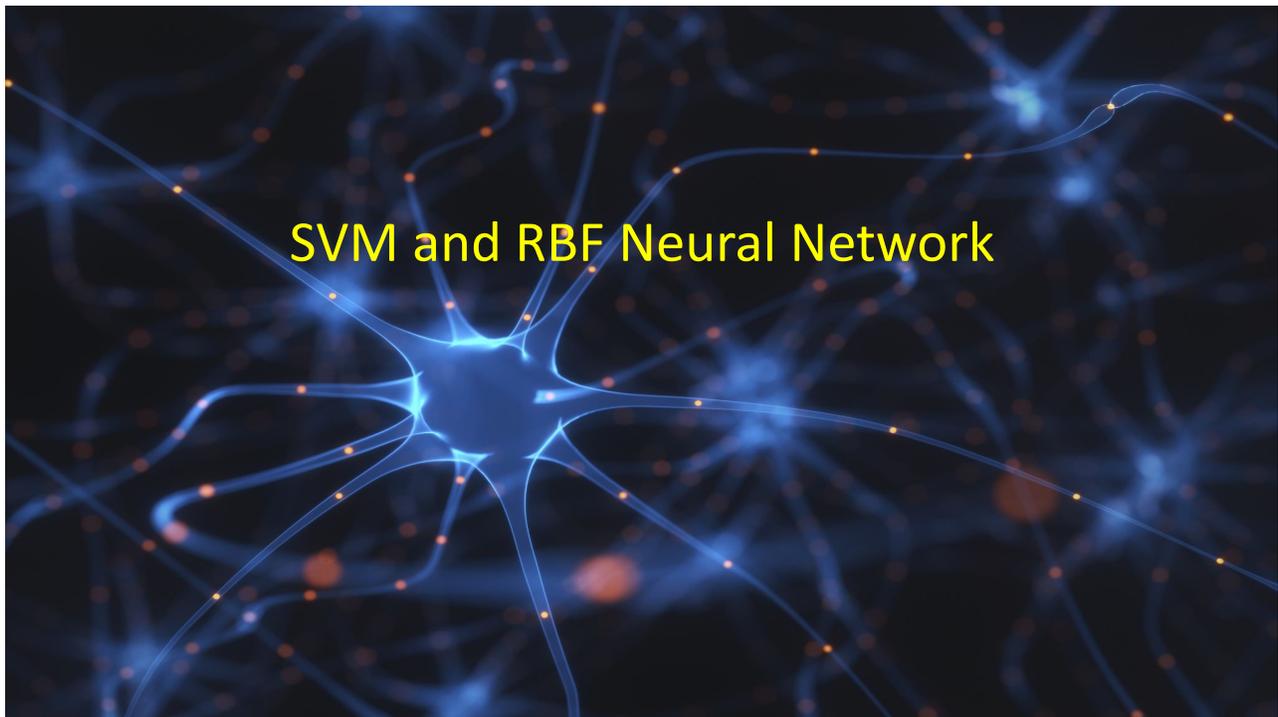
Each Region       Average of the target variable, $y_i$ in Region $R_m$       If the test feature vector is in Region $R_m$

*Compare with SVM*

$$f(x) = \text{sgn}\left( \sum_{i=1}^{n} \alpha_i y_i K(x_i, x) \right)$$

SVM and RBF Neural Network

$$RBNN: f(x) = \sum_i w_i e^{-(\frac{||x-x_i||^2}{2\sigma^2})}$$

*SVM RBF Kernel*

$$f(x) = \text{sgn}(\sum_{i=1}^n \alpha_i y_i \, e^{-(\frac{||x-x_i||^2}{2\sigma^2})})$$



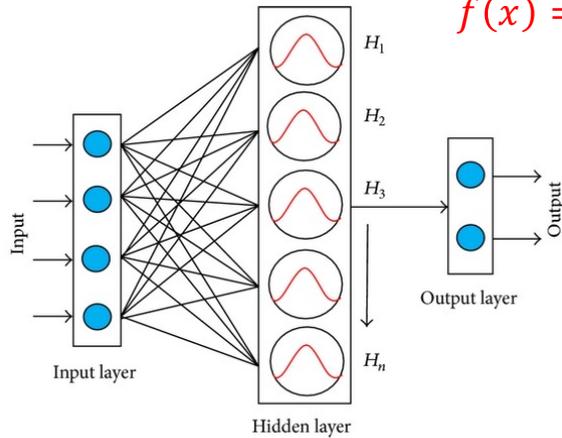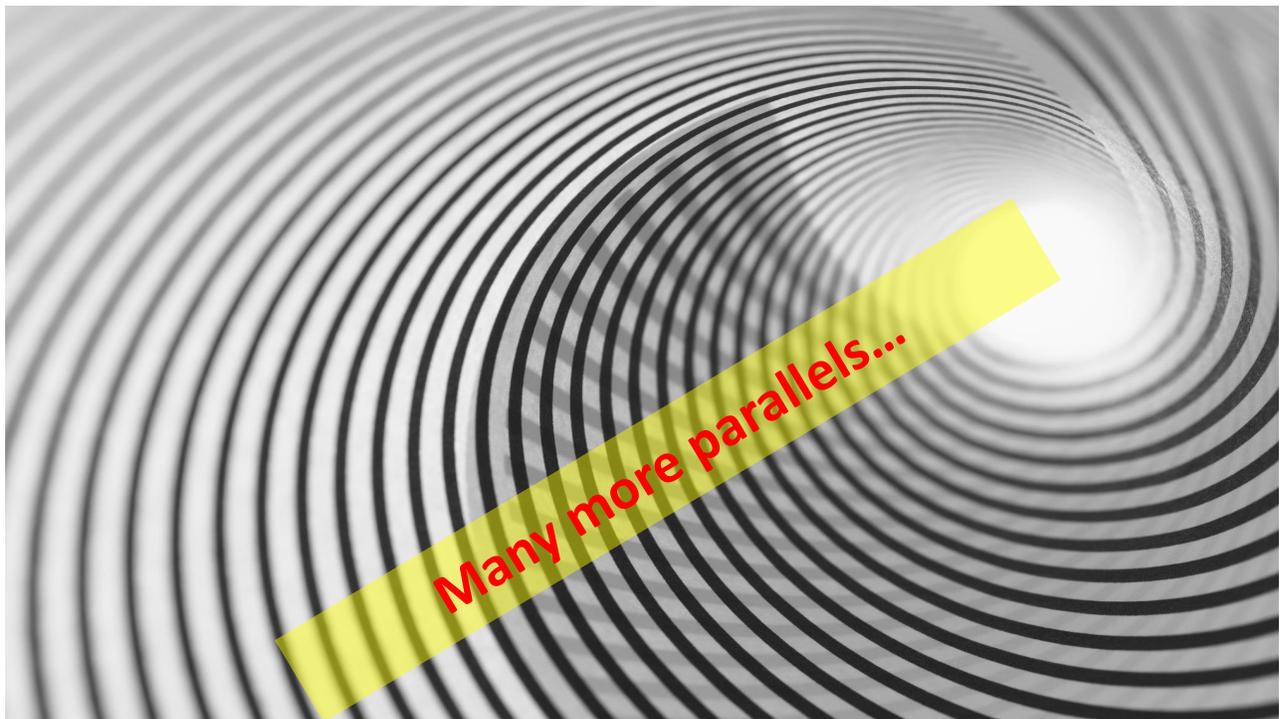*Support vectors are the hidden nodes!*

**To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.**



Many more parallels...

*https://www.sjsu.edu/people/vishnu.pendyala/*
*@vishnupendyala*

11/12/23

To cite this presentation: Pendyala, V.S. (2022) "Exploring the math in Support Vector Machines". IEEE Computer Society, Santa Clara Valley Chapter Webinar.



11/12/23