

Machine Learning

Support Vector Machines

Dr. Rizwan Ahmed Khan

Outline

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

6 Kernel Trick

- Python
- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Reference Books
Reference Books

Reference books for this Module:

- **Chapter 3:** Pattern Recognition, S. Theodoridis et al., Academic Press, 4th or latest edition.

Reference books for this Module:

- **Chapter 3:** Pattern Recognition, S. Theodoridis et al., Academic Press, 4th or latest edition.
- **Chapter 6 & 7:** Pattern Recognition and Machine Learning, Christopher M. Bishop, Springer Books, latest edition.

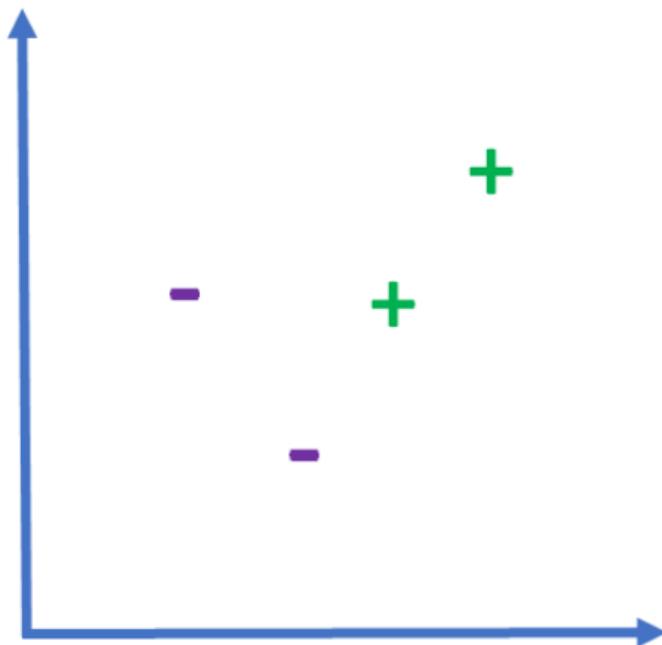
Reference Books

Reference books for this Module:

- **Chapter 3:** Pattern Recognition, S. Theodoridis et al., Academic Press, 4th or latest edition.
 - **Chapter 6 & 7:** Pattern Recognition and Machine Learning, Christopher M. Bishop, Springer Books, latest edition.
 - **Book** Support Vector Machines Succinctly, Alexandre Kowalczyk, 2017.

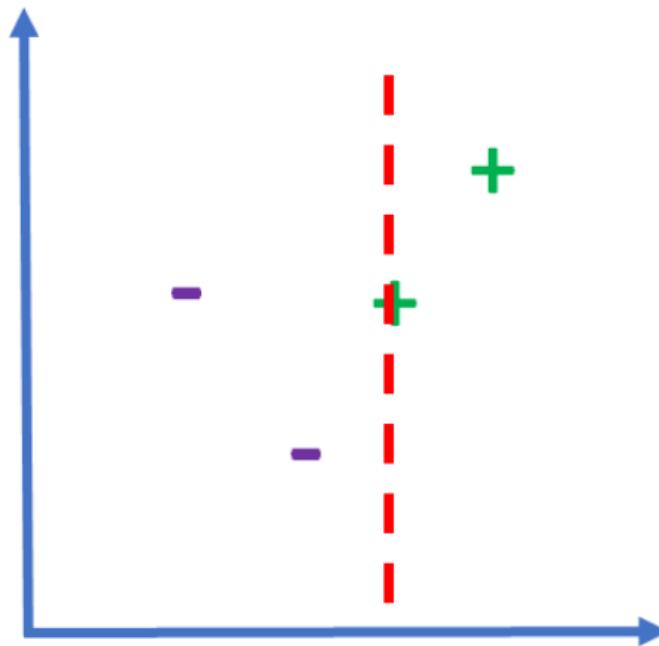
Intuition - Decision Boundary

How would you divide +ve examples from -ve examples?



Intuition - Decision Boundary

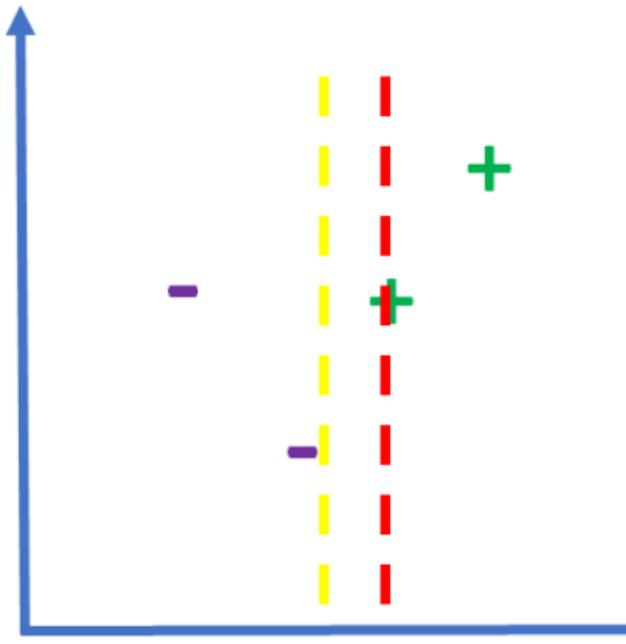
How would you divide +ve examples from -ve examples?



Seems infinite possibilities.

Intuition - Decision Boundary

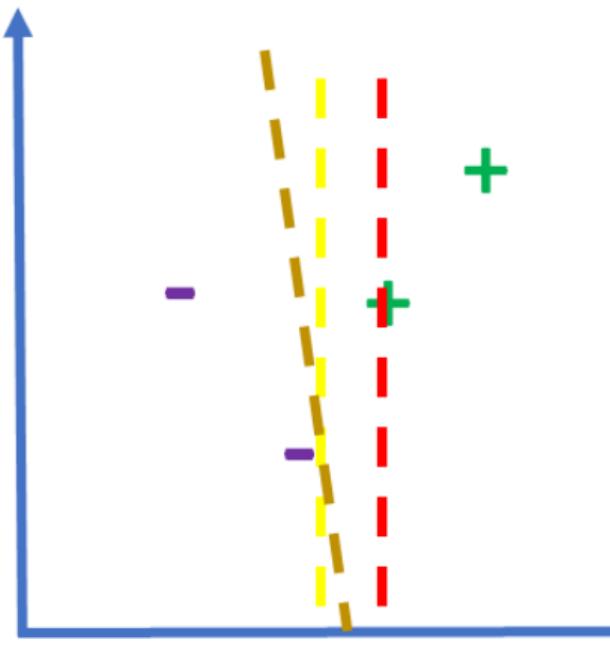
How would you divide +ve examples from -ve examples?



Seems infinite possibilities.

Intuition - Decision Boundary

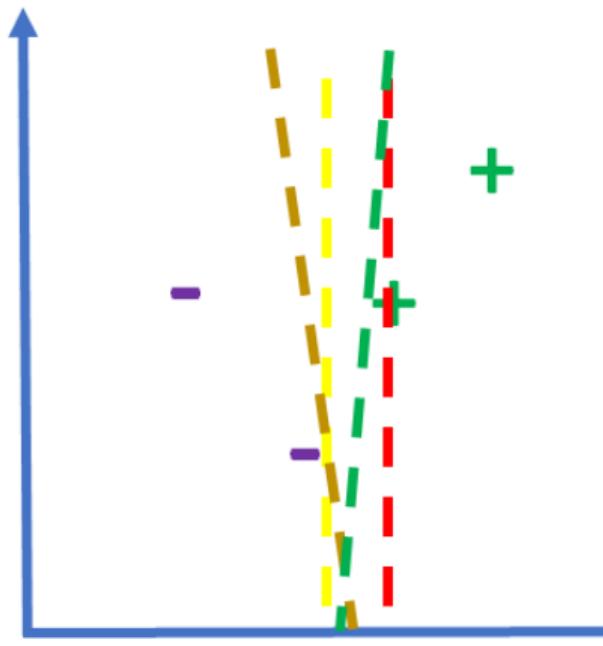
How would you divide +ve examples from -ve examples?



Seems infinite possibilities.

Intuition - Decision Boundary

How would you divide +ve examples from -ve examples?

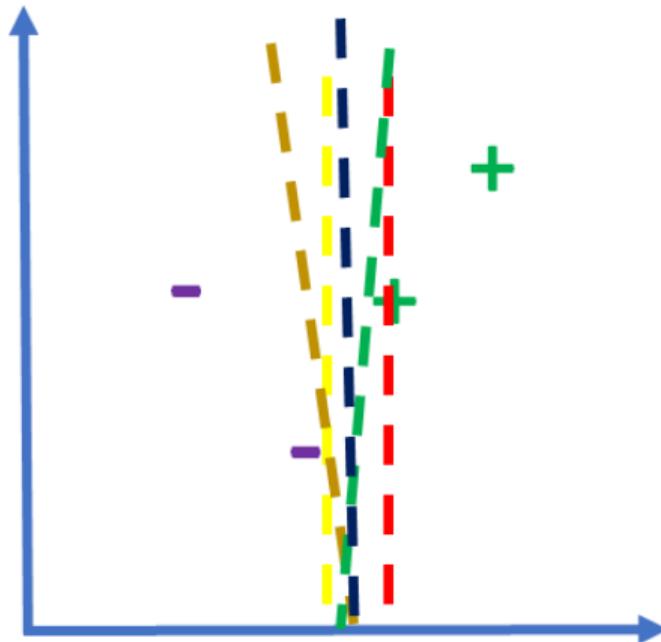


Seems infinite possibilities.

Intuition - Decision Boundary

Intuition - Decision Boundary

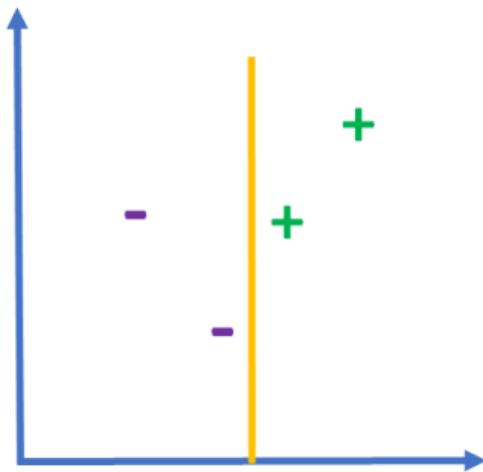
How would you divide +ve examples from -ve examples?



Seems infinite possibilities.

Intuition - Decision Boundary

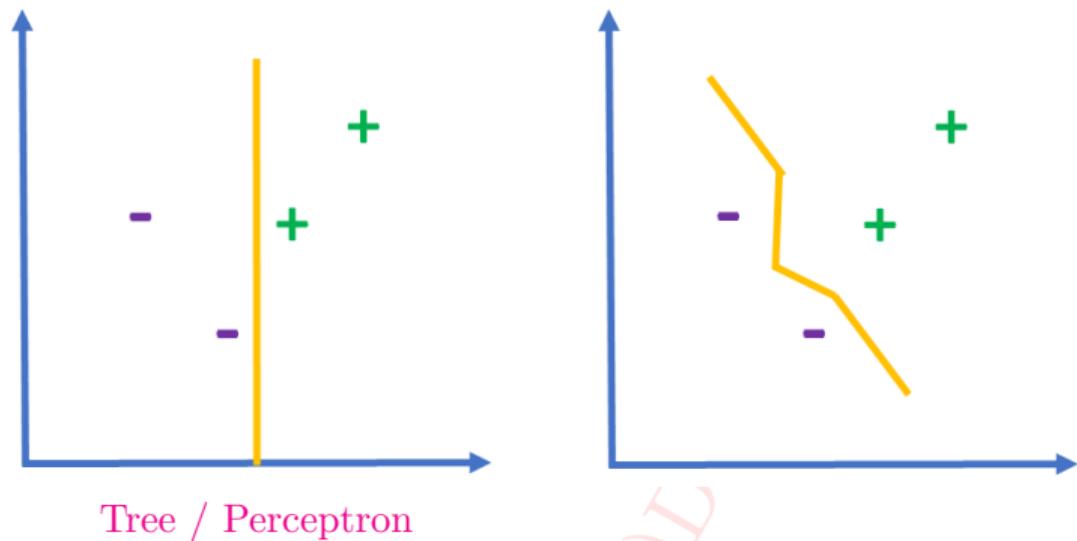
Intuition - Decision Boundary



(c)Dr. Rizwan A Khan

Intuition - Decision Boundary

Intuition - Decision Boundary

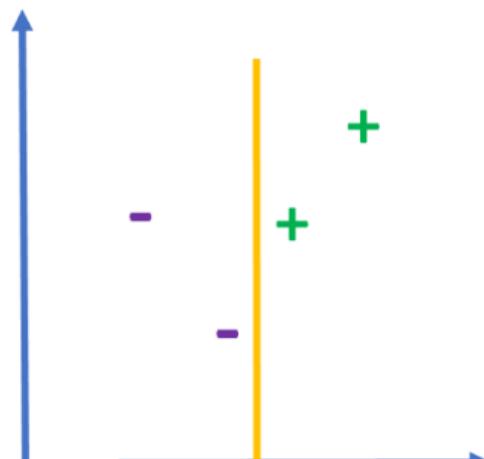


han

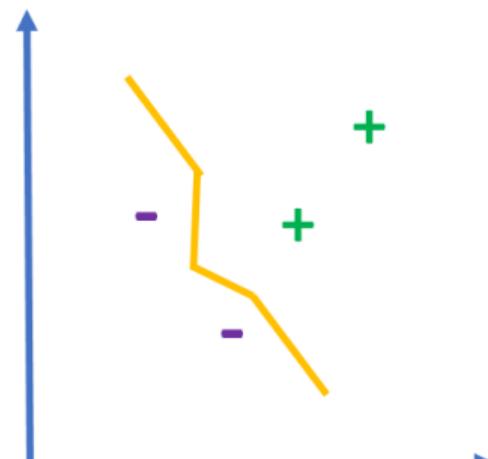
(c)L

Intuition - Decision Boundary

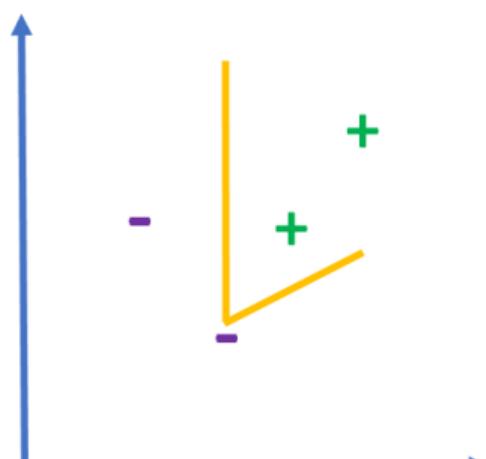
Intuition - Decision Boundary



Tree / Perceptron

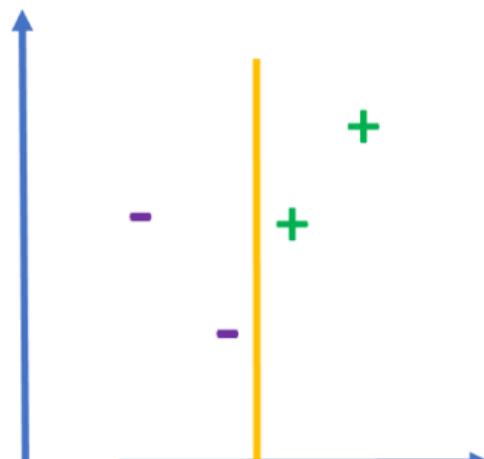


k-Nearest Neighbor

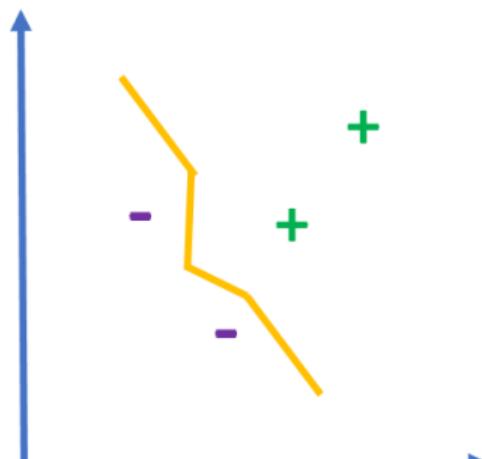
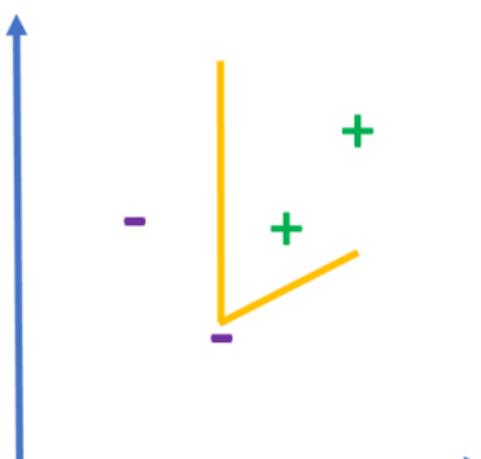


Intuition - Decision Boundary

Intuition - Decision Boundary



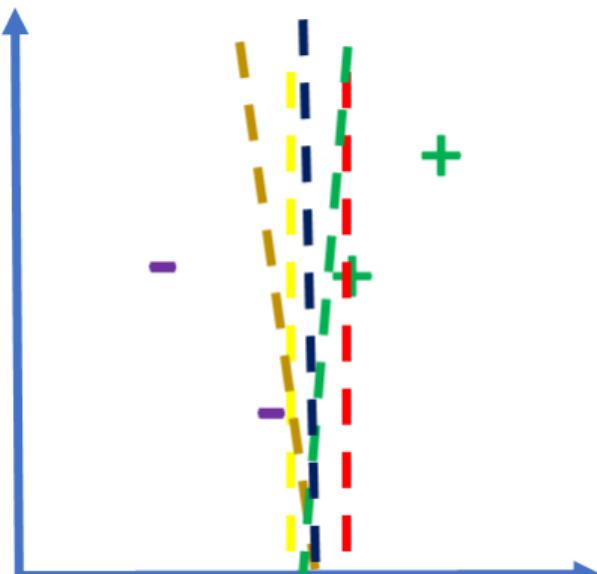
Tree / Perceptron

(c) k - Nearest Neighbor

Neural Network

Intuition - Decision Boundary

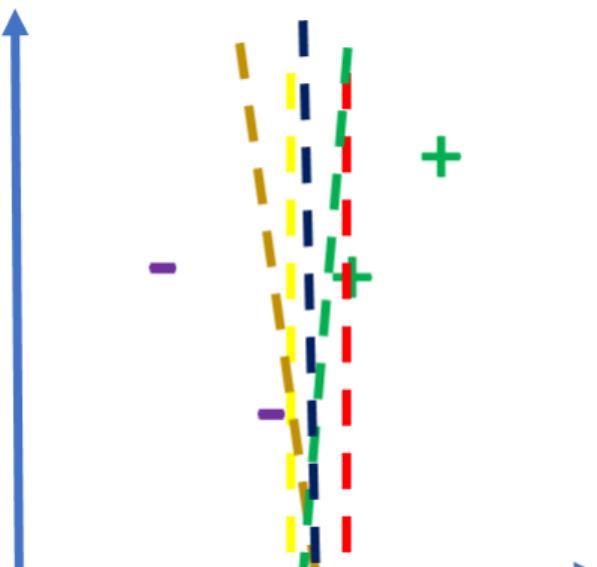
Goal of SVM



- Goal of SVM is to identify an optimal separating hyperplane which maximizes the margin between different classes of the training data.

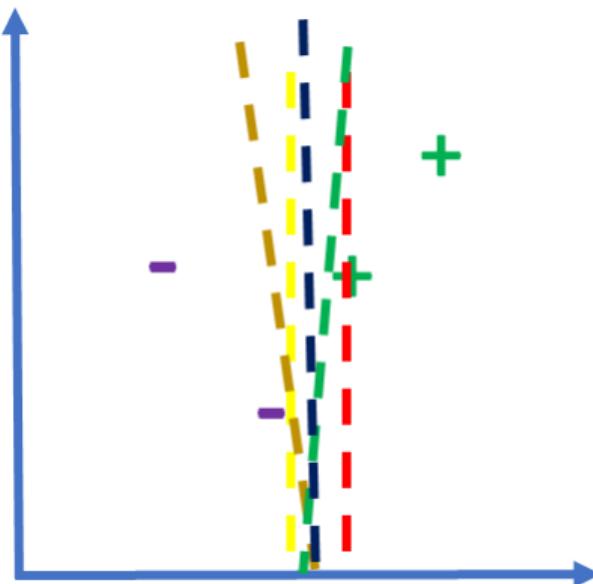
Intuition - Decision Boundary

Goal of SVM



- Goal of SVM is to identify an optimal separating hyperplane which maximizes the margin between different classes of the training data.
- SVM ^a is completely based on Mathematical Optimization problem.

Goal of SVM



- Goal of SVM is to identify an optimal separating hyperplane which maximizes the margin between different classes of the training data.
- SVM ^a is completely based on Mathematical Optimization problem.
- SVMs are linear classifiers (a line in 2 dimensions, a plane in 3 dimensions, a $n - 1$ dimensional hyperplane in n dimensions ^b).

^aCortes, C., Vapnik, V. Support-vector networks. Machine Learning 20, 273–297 (1995). <https://doi.org/10.1007/BF00994018>

^bcs276a SVM Review-Stanford University

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

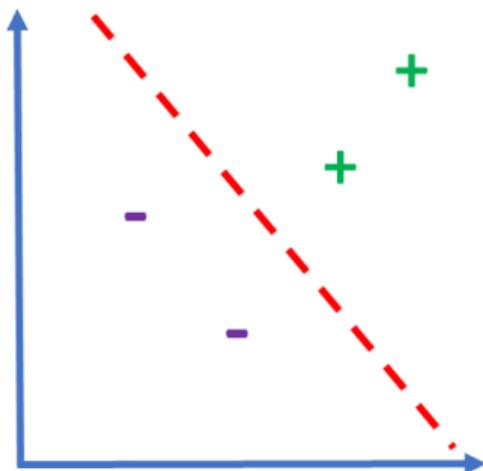
7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

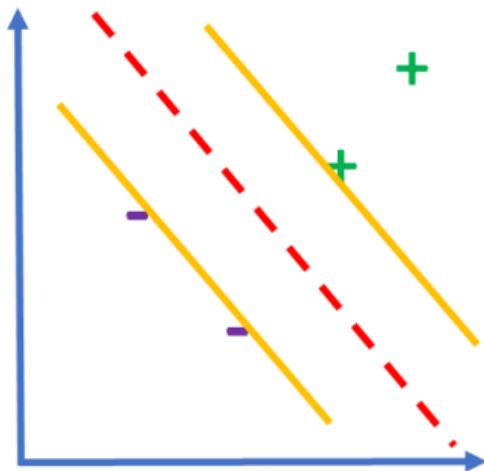
Intuition - Decision Boundary



Widest Street Approach

- Find such a line that maximizes the distance between +ve examples and -ve examples, while deciding decision boundary / surface.

Intuition - Decision Boundary

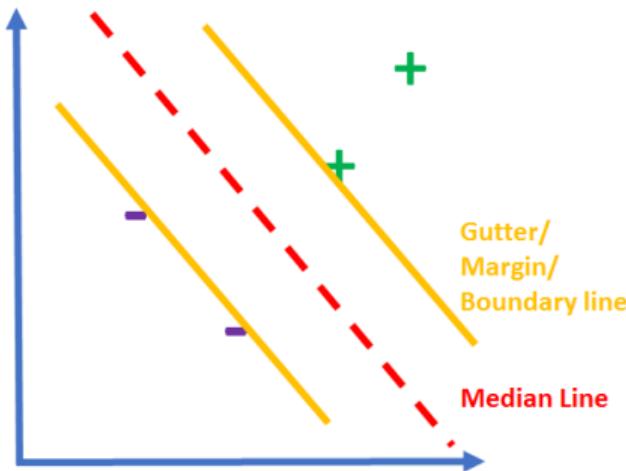


Widest Street Approach

- Find such a line that maximizes the distance between *+ve* examples and *-ve* examples, while deciding decision boundary / surface.
- What would the decision rule?

SVM - Decision Boundary

Intuition - Decision Boundary

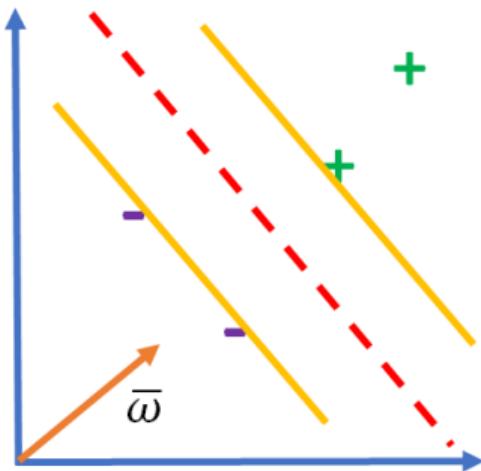


Widest Street Approach

- Find such a line that maximizes the distance between +ve examples and -ve examples, while deciding decision boundary / surface.
- What would the decision rule?

SVM - Decision Boundary

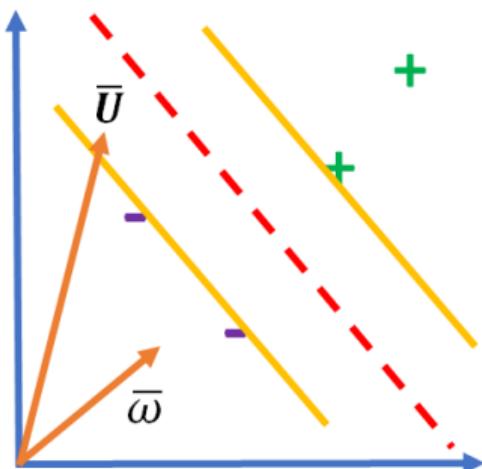
Decision Boundary



- Consider a vector \bar{W} that is perpendicular to median / or gutter. We don't know anything about its length yet.

SVM - Decision Boundary

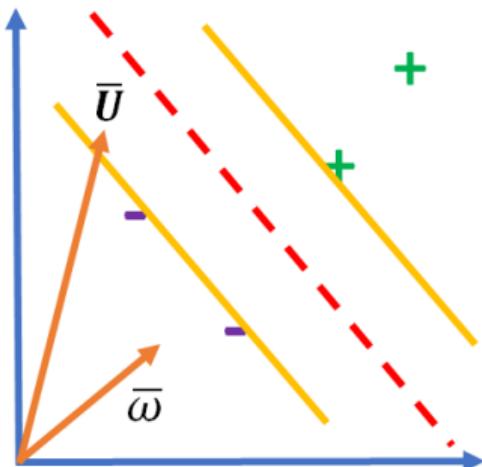
Decision Boundary



- Consider a vector \bar{W} that is perpendicular to median / or gutter. We don't know anything about its length yet.
- Consider unknown point \bar{U} and a vector points to it.

SVM - Decision Boundary

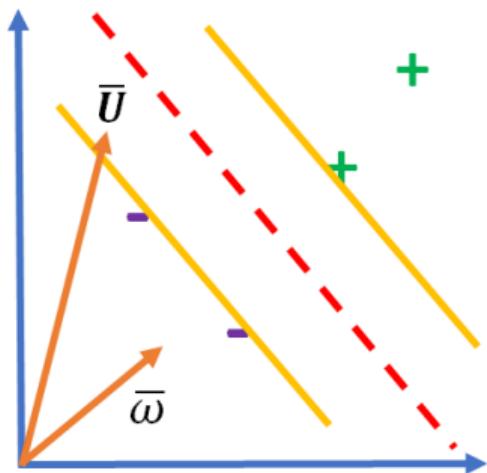
Decision Boundary



- Consider a vector \bar{W} that is perpendicular to median / or gutter. We don't know anything about its length yet.
- Consider unknown point \bar{U} and a vector points to it.
- We are interested to know whether this unknown is either right side of street or left or we want to know its label.
- What we can do, project that to perpendicular vector. The further we go we can find that its on the right side of the street.

SVM - Decision Boundary

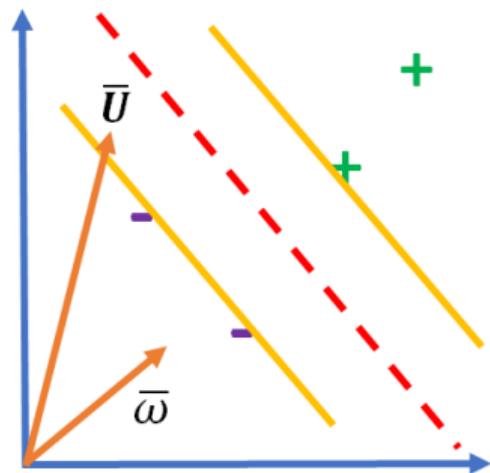
Decision Boundary



- We are interested to know whether this unknown is either right side of street or left or we want to know its label.
- What we can do, project \bar{U} to vector \bar{W} which is perpendicular median line. The further we go, we can find that its on the right side of the street.

SVM - Decision Boundary

SVM Decision Boundary

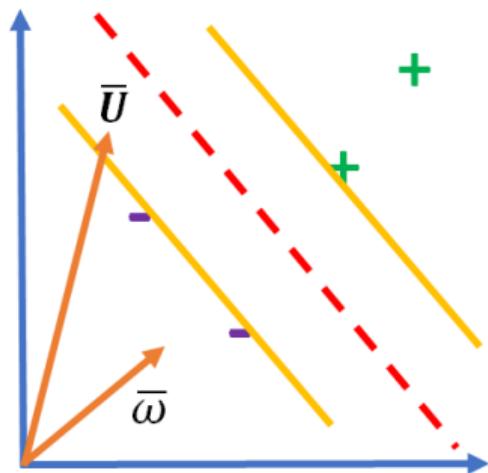


- We are interested to know whether this unknown is either right side of street or left or we want to know its label.
 - What we can do, project \bar{U} to vector \bar{W} which is perpendicular median line. The further we go, we can find that its on the right side of the street.

$$\bar{W} \cdot \bar{U} > C$$

SVM - Decision Boundary

Decision Boundary



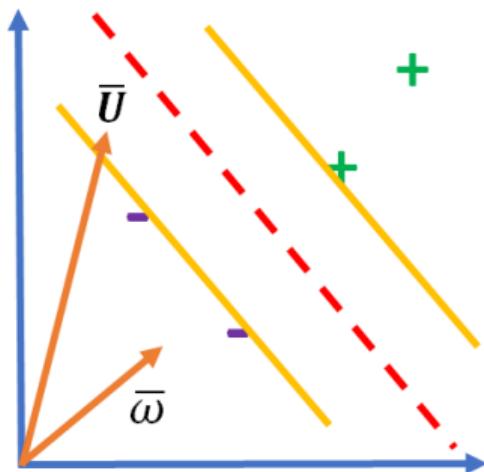
- We are interested to know whether this unknown is either right side of street or left or we want to know its label.
- What we can do, project \bar{U} to vector \bar{W} which is perpendicular median line. The further we go, we can find that its on the right side of the street.

$$\bar{W} \cdot \bar{U} \geq C$$

- Dot product is projecting onto \bar{W} . The bigger the projection is then it will cross median line and then unknown vector can be labeled as **+ve sample**.

SVM - Decision Boundary

Decision Boundary

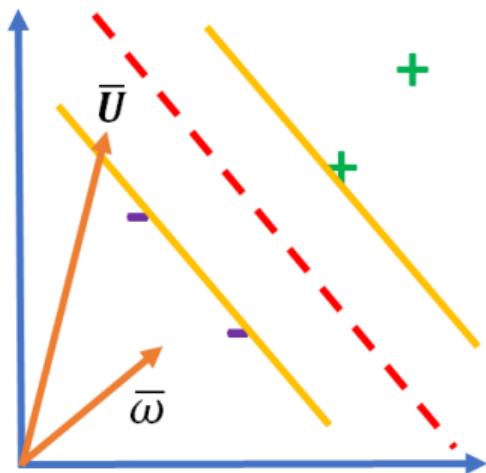


- Then, without loss of generality we can say:

(c) Dr. Rizwan Ahmed Khan

SVM - Decision Boundary

Decision Boundary



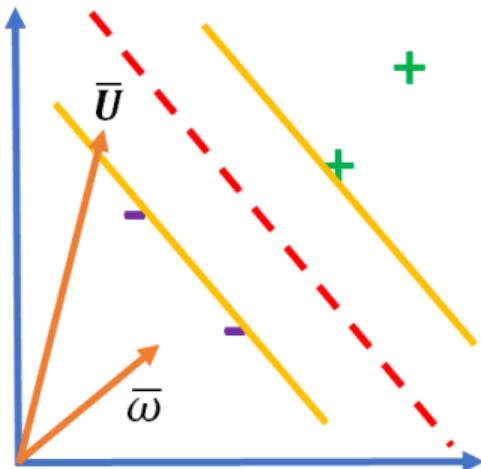
- Then, without loss of generality we can say:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve} \quad (1)$$

This is **Decision Rule**.

SVM - Decision Boundary

Decision Boundary



- Then, without loss of generality we can say:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve} \quad (1)$$

This is **Decision Rule**.

- We don't know (yet) what constant b , $C = -b$, to use and what \bar{W} to use either.

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

6 Kernel Trick

- Python
- Code
- Results
- XOR problem visualization
- Conclusion

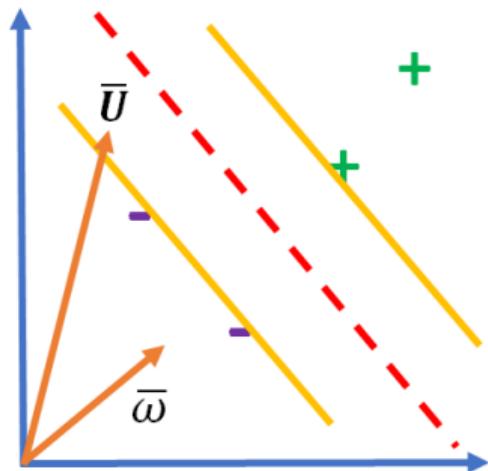
7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

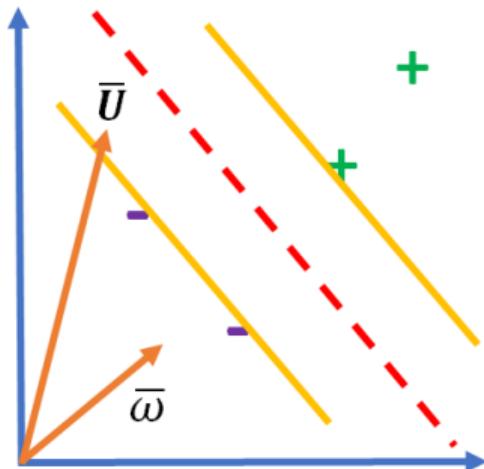
- Function Visualization
- Function Optimization

Constraints



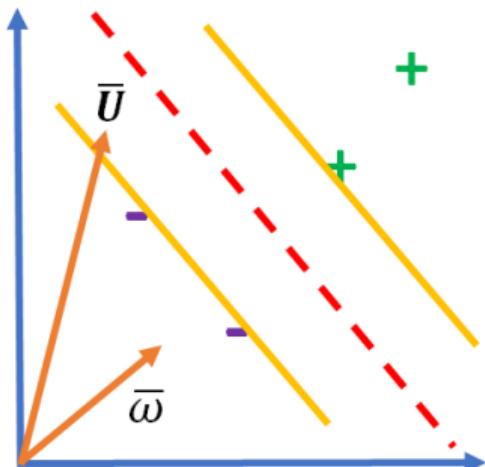
- We just know that \bar{W} needs to be perpendicular to the median line of the street.

Constraints



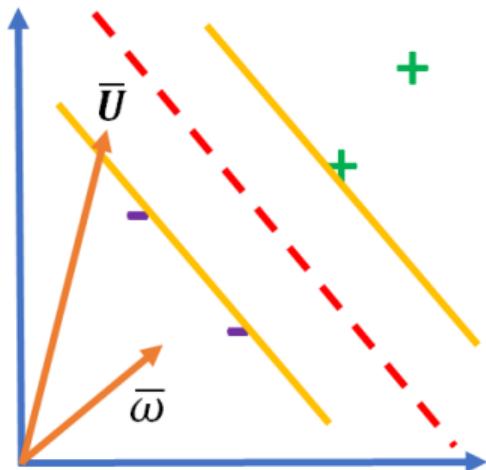
- We just know that \bar{W} needs to be perpendicular to the median line of the street.
- But then there many \bar{W}_s perpendicular to the median line of the street, any length is not fixed yet.
- What should we do?

Constraints



- We just know that \bar{W} needs to be perpendicular to the median line of the street.
- But then there many \bar{W}_s perpendicular to the median line of the street, any length is not fixed yet.
- What should we do?
- So we need to put **constraints** to find particular \bar{W} and b that maximizes width of the street (separation between $+_s$ and $-_s$).

Constraints



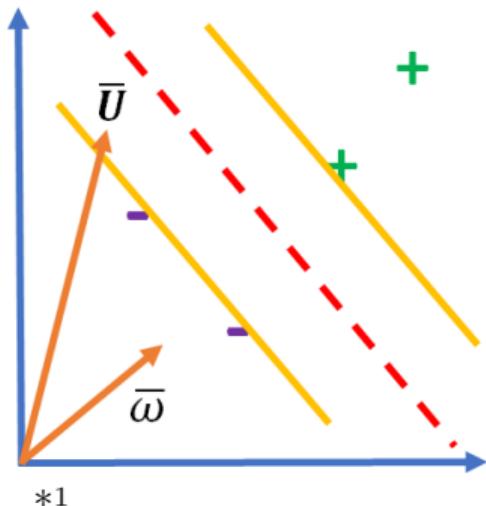
- Putting constraints to calculate \bar{W} and b .

$$\bar{W} \cdot \bar{X}_+ + b \geq 1 \quad \text{(for +ve samples)} \quad (2)$$

$$\bar{W} \cdot \bar{X}_- + b \leq -1 \quad \text{(for -ve samples)} \quad (3)$$

¹Did you see similarity with “Perceptron” decision rule?

Constraints



- Putting constraints to calculate \bar{W} and b .

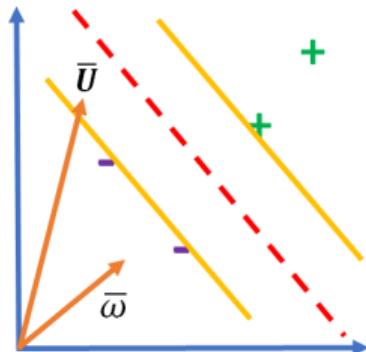
$$\bar{W} \cdot \bar{X}_+ + b \geq 1 \quad \{\text{for +ve samples}\} \quad (2)$$

$$\bar{W} \cdot \bar{X}_- + b \leq -1 \quad \{\text{for -ve samples}\} \quad (3)$$

So imposing separation of -1 to $+1$ for $-ve$ and $+ve$ samples (maximizing margin).

¹Did you see similarity with “Perceptron” decision rule?

Constraints



- Equation 2 and 3 can be written / combined as:

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1 \quad (4)$$

where:

- y_i is **+1** for *+ve* samples.
- y_i is **-1** for *-ve* samples.

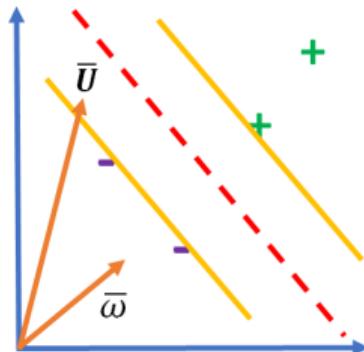
Proof:

- for *+ve* samples:

$$\begin{aligned} 1 \times (\bar{W} \cdot \bar{X}_i + b) &\geq 1 \\ \Rightarrow (\bar{W} \cdot \bar{X}_i + b) &\geq 1 \end{aligned} \quad (5)$$

Same as Eq. 2.

Constraints



- Equation 2 and 3 can be written / combined as:

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1 \quad (4)$$

where:

- y_i is **+1** for *+ve* samples.
- y_i is **-1** for *-ve* samples.

Proof:

- for *+ve* samples:

$$\begin{aligned} 1 \times (\bar{W} \cdot \bar{X}_i + b) &\geq 1 \\ \Rightarrow (\bar{W} \cdot \bar{X}_i + b) &\geq 1 \end{aligned} \quad (5)$$

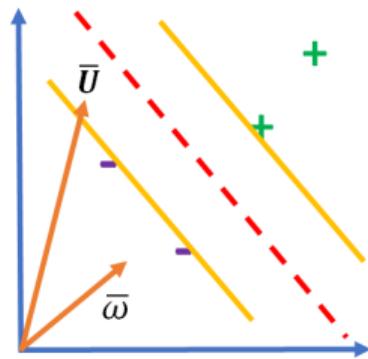
Same as Eq. 2.

- for *-ve* samples:

$$\begin{aligned} -1 \times (\bar{W} \cdot \bar{X}_i + b) &\geq 1 \\ \Rightarrow -\bar{W} \cdot \bar{X}_i - b &\geq 1 \\ \Rightarrow \bar{W} \cdot \bar{X}_i + b &\leq -1 \end{aligned} \quad (6)$$

Same as Eq. 3.

Constraints



Back to equation 4:

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1$$

where:

- y_i is +1 for +ve samples.
- y_i is -1 for -ve samples.

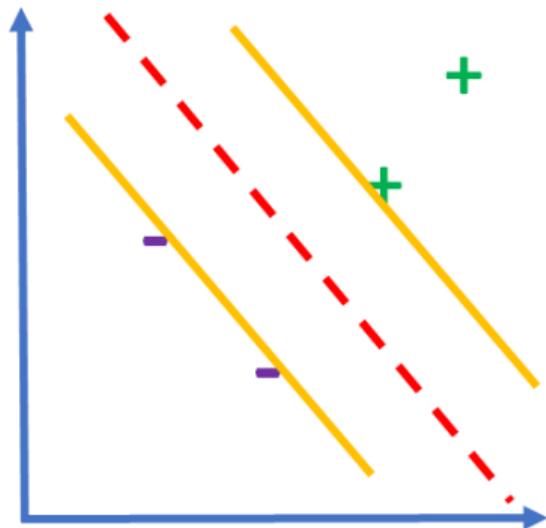
Equation 4 can be written as:

$$\begin{aligned} y_i(\bar{W} \cdot \bar{X}_i + b) &\geq 1 \\ y_i(\bar{W} \cdot \bar{X}_i + b) - 1 &\geq 0 \end{aligned} \tag{7}$$

Additional constraint:

$$\begin{aligned} y_i(\bar{W} \cdot \bar{X}_i + b) - 1 &= 0 \\ \{ \text{for samples } (X_i) \text{ in gutter or at boundary / margin} \} \end{aligned} \tag{8}$$

Samples on the boundary / gutter

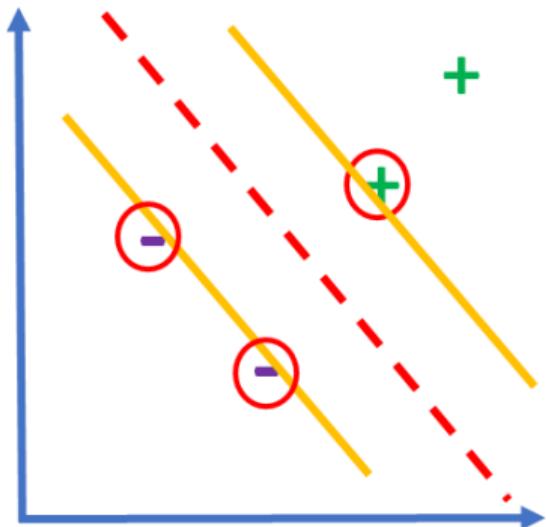


Additional constraint:

$$y_i(\bar{W} \cdot \bar{X}_i + b) - 1 = 0$$

{for samples (X_i) in gutter or at boundary / margin}

Samples on the boundary / gutter

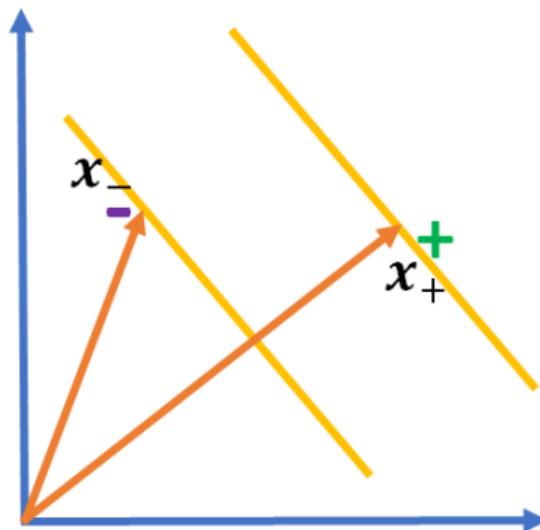


Additional constraint:

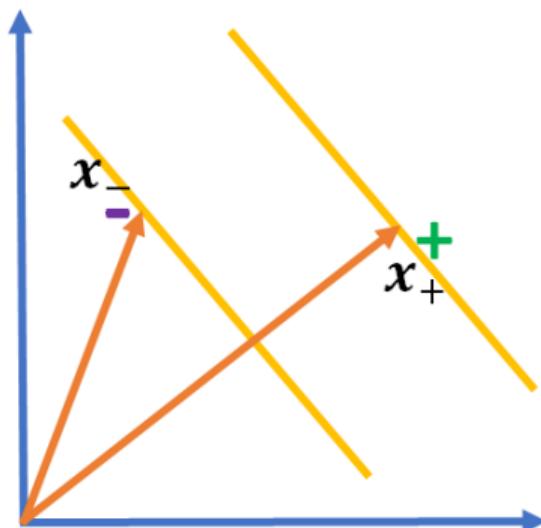
$$y_i(\bar{W} \cdot \bar{X}_i + b) - 1 = 0$$

{for samples (X_i) in gutter or at boundary / margin}

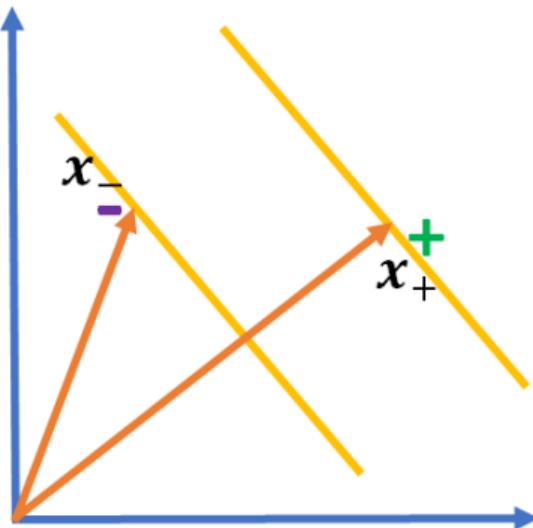
These samples are also called as **Support Vectors**.



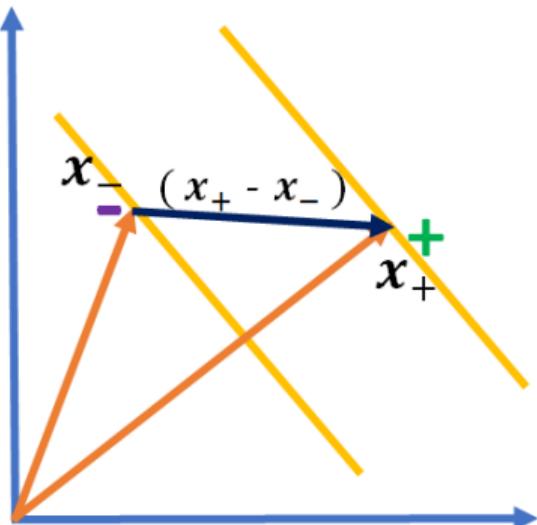
- We are trying to arrange line \bar{W} and b in a such a way that it maximizes **width of the street** (separation between $+_s$ and $-_s$).



- We are trying to arrange line \bar{W} and b in a such a way that it maximizes **width of the street** (separation between $+_s$ and $-_s$).
- Boundary lines are parallel to one another, we can pick points on these lines to define width of the street.

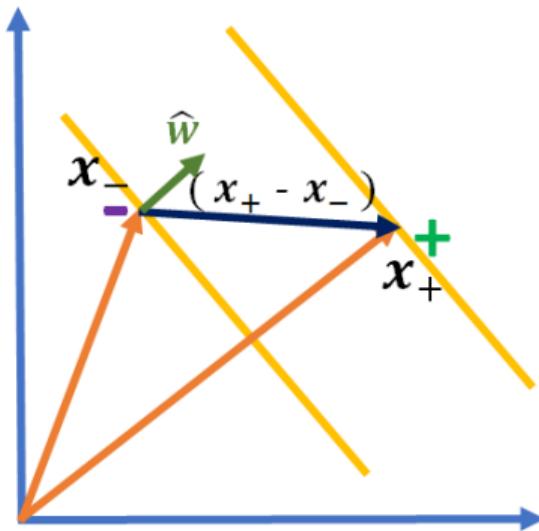


- We are trying to arrange line \bar{W} and b in a such a way that it maximizes **width of the street** (separation between $+_s$ and $-_s$).
- Boundary lines are parallel to one another, we can pick points on these lines to define width of the street.
- Width of the street is **distance b/w the gutters / boundary lines**.



- We are trying to arrange line \bar{W} and b in a such a way that it maximizes **width of the street** (separation between $+_s$ and $-_s$).
- Boundary lines are parallel to one another, we can pick points on these lines to define width of the street.
- Width of the street is **distance b/w the gutters / boundary lines**.
- Difference of two vectors can give us width of the street:

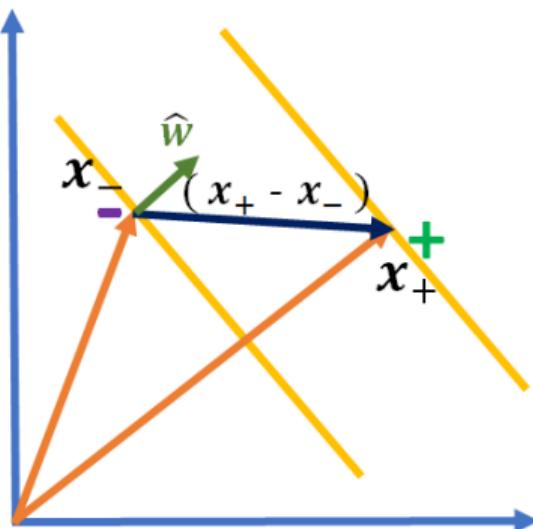
$$(\bar{X}_+ - \bar{X}_-)$$



- Width of street:

$$= (\bar{X}_+ - \bar{X}_-) \cdot \frac{\bar{W}}{\|W\|} \quad (9)$$

- Where $\frac{\bar{W}}{\|W\|}$ is a unit vector (\hat{W}) perpendicular / normal to gutter or boundary line.

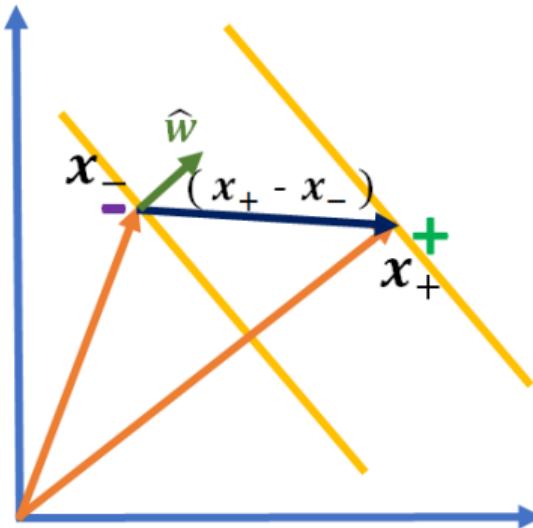


- Width of street:

$$= (\bar{X}_+ - \bar{X}_-) \cdot \frac{\bar{W}}{\|\bar{W}\|} \quad (9)$$

- Where $\frac{\bar{W}}{\|\bar{W}\|}$ is a unit vector (\hat{W}) perpendicular / normal to gutter or boundary line.

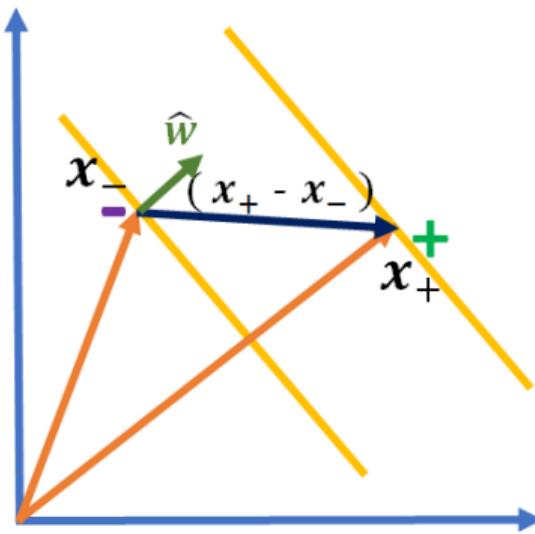
- In other words, projection of difference vector on to unit vector (\hat{W}) will be width of the street (difference in the direction of \bar{W} vector).



- Width of street:

$$= (\bar{X}_+ - \bar{X}_-) \cdot \frac{\bar{W}}{\|W\|} \quad (9)$$

- Where $\frac{\bar{W}}{\|W\|}$ is a unit vector (\hat{W}) perpendicular / normal to gutter or boundary line.
 - In other words, projection of difference vector on to unit vector (\hat{W}) will be width of the street (difference in the direction of \bar{W} vector).
 - It's a dot product, so its scalar, width of the street.



Width of street:

$$= (\bar{X}_+ - \bar{X}_-) \cdot \frac{\bar{W}}{\|W\|}$$

From Equation 8 we know that, samples in a gutter (enforcing constraint) =

$$y_i(\bar{W} \cdot \bar{X}_i + b) - 1 = 0$$

So,

- for +ve sample:

$$\bar{W} \cdot \bar{X}_i = 1 - b \quad (10)$$

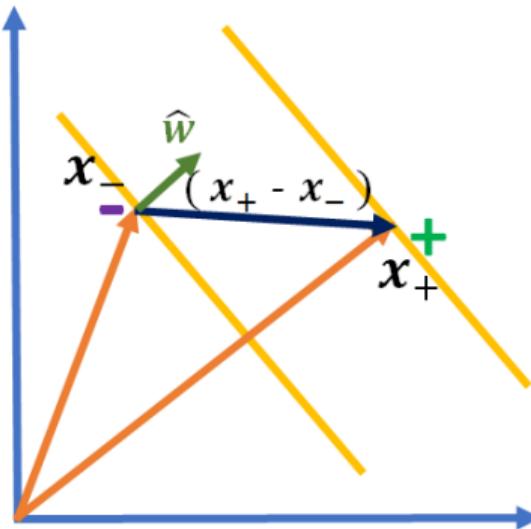
- for -ve sample:

$$-\bar{W} \cdot \bar{X}_i - b - 1 = 0$$

$$-\bar{W} \cdot \bar{X}_i \equiv 1 + b \quad (11)$$

Distance b/w boundary lines

Distance b/w boundary lines / margin / gutters

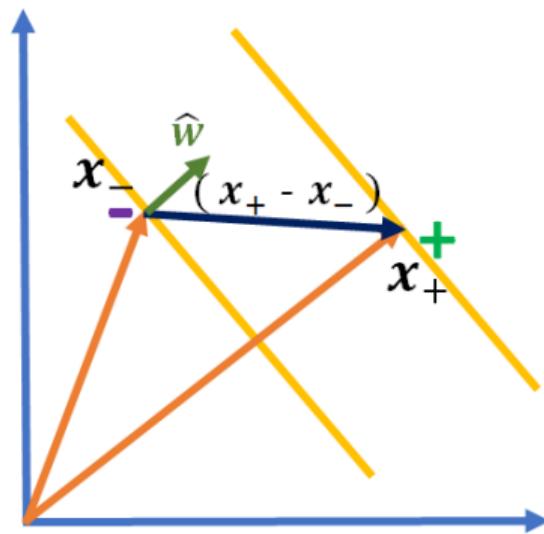


Width of street:

$$\begin{aligned} & (\bar{X}_+ - \bar{X}_-) \cdot \frac{\bar{W}}{\|W\|} \\ & \bar{W} \cdot \bar{X}_+ - \bar{W} \cdot \bar{X}_- \cdot \frac{1}{\|W\|} \end{aligned} \quad (12)$$

Putting back values from Equations 10 and 11 into Equation 12.

$$\text{Width} = \frac{2}{\|W\|} \quad (13)$$



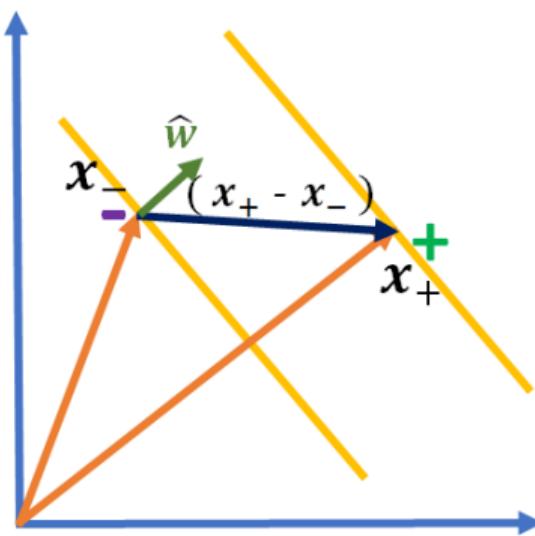
Width of street:

$$\text{Width} = \frac{2}{\|W\|}$$

- SVM tries to **maximize this width**, to have maximum possible separation between samples of different classes.

$$\text{Width} = \text{Max} \frac{2}{\|W\|} \implies \text{Min} \|W\| \implies \text{Min} \frac{1}{2} \|W\|^2$$

Width of street:



$$\text{Width} = \frac{2}{\|W\|}$$

- SVM tries to **maximize this width**, to have maximum possible separation between samples of different classes.

$$\text{Width} = \text{Max} \frac{2}{\|W\|} \implies \text{Min} \|W\| \implies \text{Min} \frac{1}{2} \|W\|^2$$

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2 \quad (14)$$

Stages in the development:

- ① Decision Rule:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve}$$

- ② Projection of W on U and put a constraint then it should be $\geq +1$ for +ve samples and ≤ -1 for -ve samples.

$$\bar{W} \cdot \bar{X}_+ + b \geq 1 \quad \{\text{for +ve samples}\}$$

$$\bar{W} \cdot \bar{X}_- + b \leq -1 \quad \{\text{for -ve samples}\}$$

- ③ Additional constraint for samples in gutter

$$y_i(\bar{W} \cdot \bar{X}_i + b) - 1 = 0$$

- ④ Then we discovered we wish to maximize / minimize this expression:

$$\text{Width} = \text{Min} \frac{1}{2} ||W||^2$$

- We have now transformed the problem into a form that can be efficiently solved.

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2$$

- The above is an optimization problem with a convex quadratic objective and some constraints. Its solution gives us the **optimal margin classifier**.

- We have now transformed the problem into a form that can be efficiently solved.

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2$$

- The above is an optimization problem with a convex quadratic objective and some constraints. Its solution gives us the **optimal margin classifier**.
- How can we go forward (**QUIZ**):

- We have now transformed the problem into a form that can be efficiently solved.

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2$$

- The above is an optimization problem with a convex quadratic objective and some constraints. Its solution gives us the **optimal margin classifier**.
- How can we go forward (**QUIZ**):
 - ① Laplace

- We have now transformed the problem into a form that can be efficiently solved.

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2$$

- The above is an optimization problem with a convex quadratic objective and some constraints. Its solution gives us the **optimal margin classifier**.
- How can we go forward (**QUIZ**):
 - ① Laplace
 - ② Legendre

- We have now transformed the problem into a form that can be efficiently solved.

$$\text{Width} = \text{Min} \frac{1}{2} \|W\|^2$$

- The above is an optimization problem with a convex quadratic objective and some constraints. Its solution gives us the **optimal margin classifier**.
- How can we go forward (**QUIZ**):
 - ① Laplace
 - ② Legendre
 - ③ Lagrange

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Lagrange Multiplier

- Lagrange multipliers provides a way for **finding extremum of a function** subject to equality constraints i.e., subject to the condition that one or more equations have to be satisfied exactly by the chosen values of the variables.
- The great advantage of this method is that it allows the **optimization to be solved without explicit parameterization in terms of the constraints**.
- Method can be summarized as follows: in order to find the stationary points of a function $f(x)$ subjected to the equality constraint $g(x) = 0$, form the Lagrangian function:

$$L(x, \lambda) = f(x) - \lambda g(x) \quad (15)$$

where λ = Lagrange multiplier

2

²Refer Section 8 to see discussion on intuition of Lagrange multiplier

Lagrange Multiplier

- Taking Equation 15 and writing function that we are trying to find extremum.

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i (\text{write down constraints})$$

Lagrange Multiplier

- Taking **Equation 15** and writing function that we are trying to find extremum.

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i (\text{write down constraints})$$

- Constraint is given in **Equation 8**

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1] \quad (16)$$

where α_i = Lagrange multiplier

α_i will be non-zero for vectors connected with samples in gutter, otherwise it will be zero.

Find Extremum

- What needs to be done to find extremum of **Equation 16** ?

$$L = \frac{1}{2} ||W||^2 - \sum \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

Find Extremum

- What needs to be done to find extremum of **Equation 16** ?

$$L = \frac{1}{2}||W||^2 - \sum \alpha_i [y_i(\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative and set it equal to zero.

Find Extremum

- What needs to be done to find extremum of **Equation 16** ?

$$L = \frac{1}{2}||W||^2 - \sum \alpha_i [y_i(\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative and set it equal to zero.
- **Take derivative w.r.t. “W”:**

Find Extremum

- What needs to be done to find extremum of Equation 16 ?

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i(\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative and set it equal to zero.

- Take derivative w.r.t. “W”:

$$\frac{\partial L}{\partial W} = \bar{W} - \sum_i \alpha_i y_i \bar{X}_i = 0 \implies W = \sum_i \alpha_i y_i \bar{X}_i \quad (17)$$

Where α_i is a scalar, y_i is +1 or -1 and X_i is sample vector. **What this equation signifies?**

Find Extremum

- What needs to be done to find extremum of **Equation 16** ?

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative and set it equal to zero.

- **Take derivative w.r.t. “W”:**

$$\frac{\partial L}{\partial W} = \bar{W} - \sum_i \alpha_i y_i \bar{X}_i = 0 \implies W = \sum_i \alpha_i y_i X_i \quad (17)$$

Where α_i is a scalar, y_i is +1 or -1 and X_i is sample vector. **What this equation signifies?**

What this equation signifies

Decision vector W is linear sum of **some** samples. Some, in the sense that α_i will be non-zero for few vectors (connected with samples in gutter).

Find Extremum

- Back to **Equation 16**. Any other variable that may vary?

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

Find Extremum

- Back to **Equation 16**. Any other variable that may vary?

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i(\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative w.r.t. “b”:

Find Extremum

- Back to **Equation 16**. Any other variable that may vary?

$$L = \frac{1}{2} \|W\|^2 - \sum \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

- Take derivative w.r.t. “b”:

$$\frac{\partial L}{\partial b} = - \sum_i \alpha_i y_i = 0 \implies \sum_i \alpha_i y_i = 0 \quad (18)$$

Find Extremum

Plug back value of W from Equation 17 to Equation 16.

$$W = \sum_i \alpha_i y_i \bar{X}_i$$

$$L = \frac{1}{2} \|W\|^2 - \sum_i \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

Find Extremum

Plug back value of W from [Equation 17](#) to [Equation 16](#).

$$W = \sum_i \alpha_i y_i \bar{X}_i$$

$$L = \frac{1}{2} \|W\|^2 - \sum_i \alpha_i [y_i (\bar{W} \cdot \bar{X}_i + b) - 1]$$

$$\begin{aligned} L &= \frac{1}{2} \left(\sum_i \alpha_i y_i \bar{X}_i \right) \cdot \left(\sum_j \alpha_j y_j \bar{X}_j \right) \\ &\quad - \left(\sum_i \alpha_i y_i \bar{X}_i \right) \cdot \left(\sum_j \alpha_j y_j \bar{X}_j \right) \\ &\quad - \sum_i \alpha_i y_i b + \sum_i \alpha_i \end{aligned} \tag{19}$$

Find Extremum

Term shown in red in **Equation 19** = 0 , refer **Equation 18**

Now arrange and re-write **Lagrangian Equation 19**

Find Extremum

Term shown in red in **Equation 19** = 0 , refer **Equation 18**

Now arrange and re-write **Lagrangian Equation 19**

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j \quad (20)$$

Find Extremum

Term shown in red in **Equation 19** = 0 , refer **Equation 18**

Now arrange and re-write **Lagrangian Equation 19**

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j \quad (20)$$

What this equation signifies

This optimization / finding extremum of a function depends only on dot products of pairs of samples (**dual problem**).

Decision Rule

Recall **Equation 1**, related to decision rule:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve}$$

Now, we can update this Equation, with derived value of \bar{W} , refer **Equation 17**:

Decision Rule

Recall **Equation 1**, related to decision rule:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve}$$

Now, we can update this Equation, with derived value of \bar{W} , refer **Equation 17**:

$$\sum_i \alpha_i y_i \bar{X}_i \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve} \quad (21)$$

Decision Rule

Recall **Equation 1**, related to decision rule:

$$\bar{W} \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve}$$

Now, we can update this Equation, with derived value of \bar{W} , refer **Equation 17**:

$$\sum_i \alpha_i y_i \bar{X}_i \cdot \bar{U} + b \geq 0 \quad \text{THEN +ve} \quad (21)$$

What this equation signifies?

Decision rule depends again only on dot product of unknown vector and sample vector.

an

Summary

- Recall **Equation 13**, Width = $\frac{2}{\|\bar{W}\|}$, while satisfying constraints, given by: classify training examples correctly $y_i(\bar{W} \cdot \bar{X}_i + b) - 1 \geq 0, \forall i$.

an

Summary

- Recall **Equation 13**, Width = $\text{Max}_{\frac{2}{||W||}}$, while satisfying constraints, given by: classify training examples correctly $y_i(\bar{W} \cdot \bar{X}_i + b) - 1 \geq 0, \forall i$.
- Then, we transformed above equation into equivalent problem (which is easier to solve), **Equation 14**, Width = $\text{Min}_{\frac{1}{2}} ||W||^2$.
- Why its easier?

an

Summary

- Recall **Equation 13**, Width = Max $\frac{2}{\|\bar{W}\|}$, while satisfying constraints, given by: classify training examples correctly $y_i(\bar{W} \cdot \bar{X}_i + b) - 1 \geq 0, \forall i$.
- Then, we transformed above equation into equivalent problem (which is easier to solve), **Equation 14**, Width = Min $\frac{1}{2}||W||^2$.
- Why its easier?
- Actually, this is easier to solve as when we have optimization problem in the form given above while satisfying constraints, is called **quadratic programming (QP)** problem. QP is well known field and it's solution is easier to find.

an

Summary

- Recall **Equation 13**, Width = Max $\frac{2}{\|\bar{W}\|}$, while satisfying constraints, given by: classify training examples correctly $y_i(\bar{W} \cdot \bar{X}_i + b) - 1 \geq 0, \forall i$.
- Then, we transformed above equation into equivalent problem (which is easier to solve), **Equation 14**, Width = Min $\frac{1}{2}||W||^2$.
- Why its easier?
- Actually, this is easier to solve as when we have optimization problem in the form given above while satisfying constraints, is called **quadratic programming (QP)** problem. QP is well known field and it's solution is easier to find.
- Optimization problems of this form have convex function and thus unique solution is always guaranteed.

an

Summary

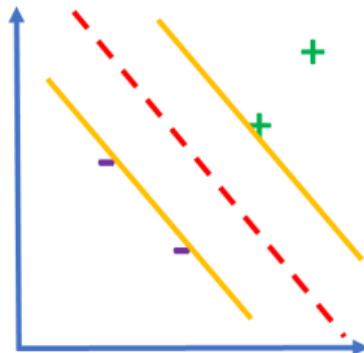
- Recall **Equation 13**, Width = Max $\frac{2}{\|\bar{W}\|}$, while satisfying constraints, given by: classify training examples correctly $y_i(\bar{W} \cdot \bar{X}_i + b) - 1 \geq 0, \forall i$.
- Then, we transformed above equation into equivalent problem (which is easier to solve), **Equation 14**, Width = Min $\frac{1}{2}||W||^2$.
- Why its easier?
- Actually, this is easier to solve as when we have optimization problem in the form given above while satisfying constraints, is called **quadratic programming (QP)** problem. QP is well known field and it's solution is easier to find.
- Optimization problems of this form have convex function and thus unique solution is always guaranteed.
- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$

Summary

- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$
While, $W = \sum_i \alpha_i y_i \bar{X}_i$ and
 $\sum_i \alpha_i y_i = 0$

Summary

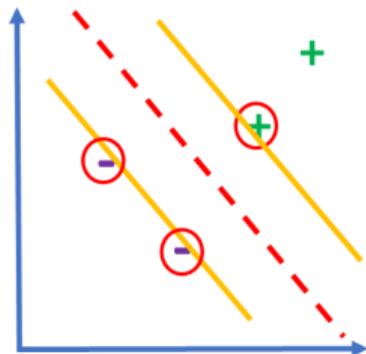
- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$
While, $W = \sum_i \alpha_i y_i \bar{X}_i$ and
 $\sum_i \alpha_i y_i = 0$



It turns out most of α_i are zeros, which implies that only few vectors (with non-zero α_i) matters in finding solution / decision boundary while most of vectors do not. Thus, building a **machine** with few **support vectors** (with non-zero α_i).

Summary

- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$
While, $W = \sum_i \alpha_i y_i \bar{X}_i$ and
 $\sum_i \alpha_i y_i = 0$



It turns out most of α_i are zeros, which implies that only few vectors (with non-zero α_i) matters in finding solution / decision boundary while most of vectors do not. Thus, building a **machine** with few **support vectors** (with non-zero α_i).

han

Summary

- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$
While, $\bar{W} = \sum_i \alpha_i y_i \bar{X}_i$ and
 $\sum_i \alpha_i y_i = 0$
- This optimization / finding extremum of a function depends only on dot products of pairs of samples , $(\bar{X}_i^T \cdot \bar{X}_j)$. Note^a.
- Analyzing : $\bar{X}_i^T \cdot \bar{X}_j$, What it actually means?

han

Summary

- Quadratic programming problem form: $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$
While, $\bar{W} = \sum_i \alpha_i y_i \bar{X}_i$ and
 $\sum_i \alpha_i y_i = 0$
- This optimization / finding extremum of a function depends only on dot products of pairs of samples , $(\bar{X}_i^T \cdot \bar{X}_j)$. Note^a.
- Analyzing : $\bar{X}_i^T \cdot \bar{X}_j$, What it actually means?
 - Its a dot product, projection of one on another.
 - It is a **measure of similarity** between two non-zero vectors. If vectors are orthogonal value will be zero, and if vector in opposite direction value will be $-ve$.

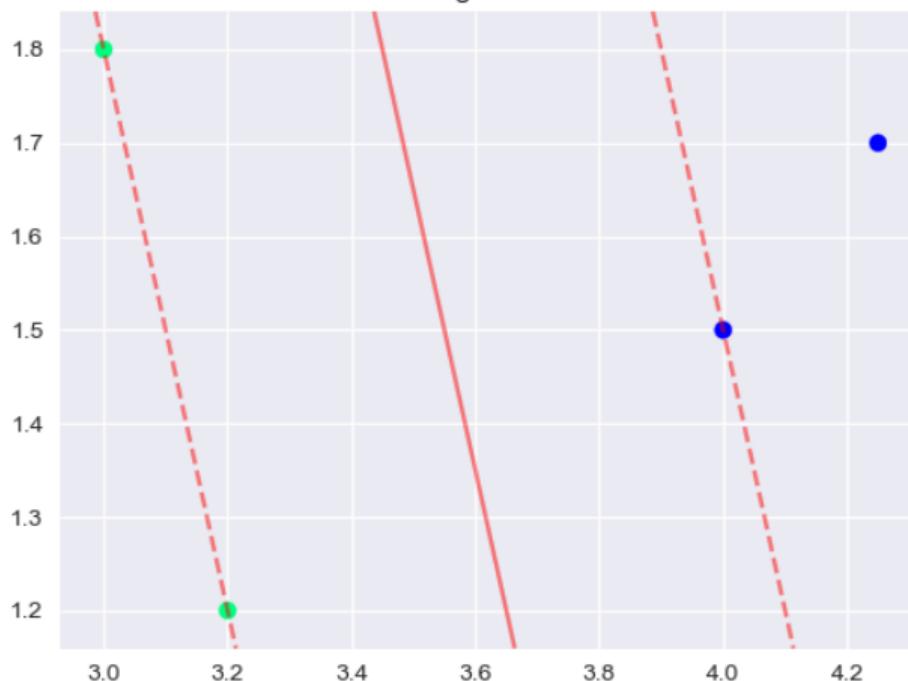
^aTranspose is used to make matrix dimensions compatible

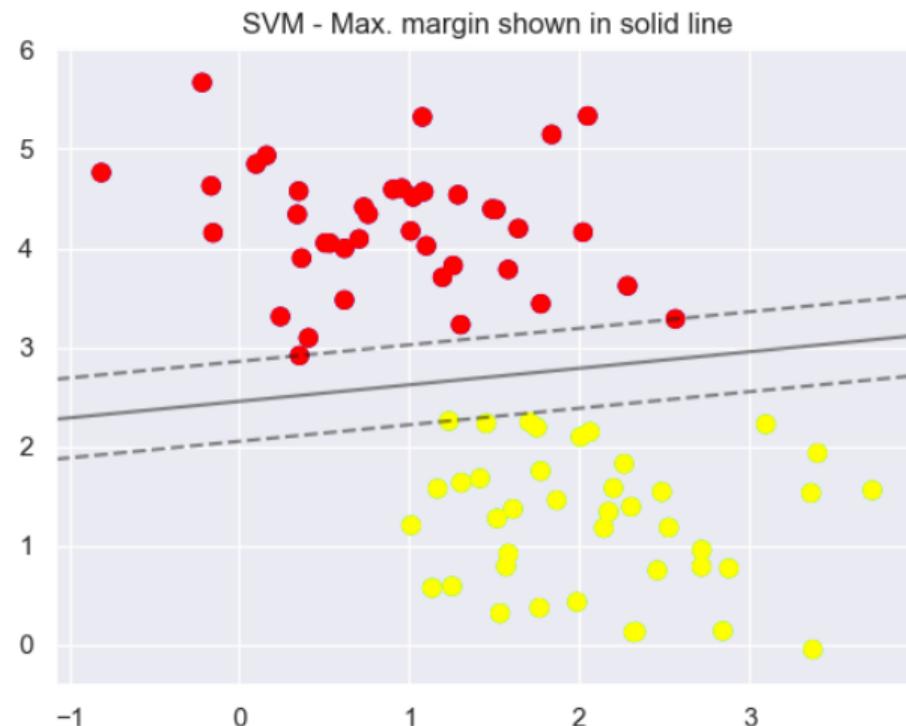
Ocular proof

Ocular proof



SVM - Max. margin shown in solid line



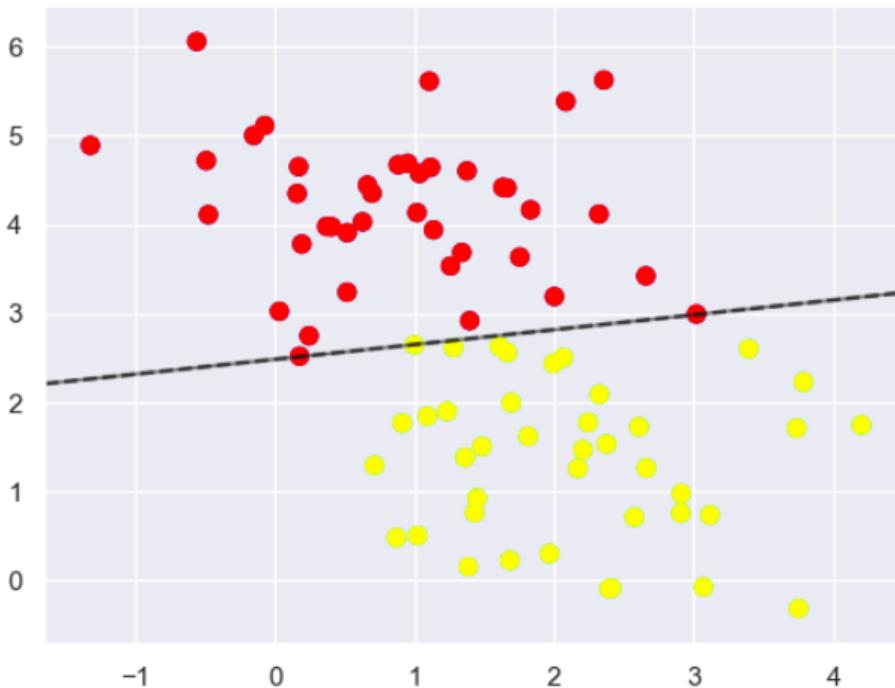


Ocular proof

Ocular proof

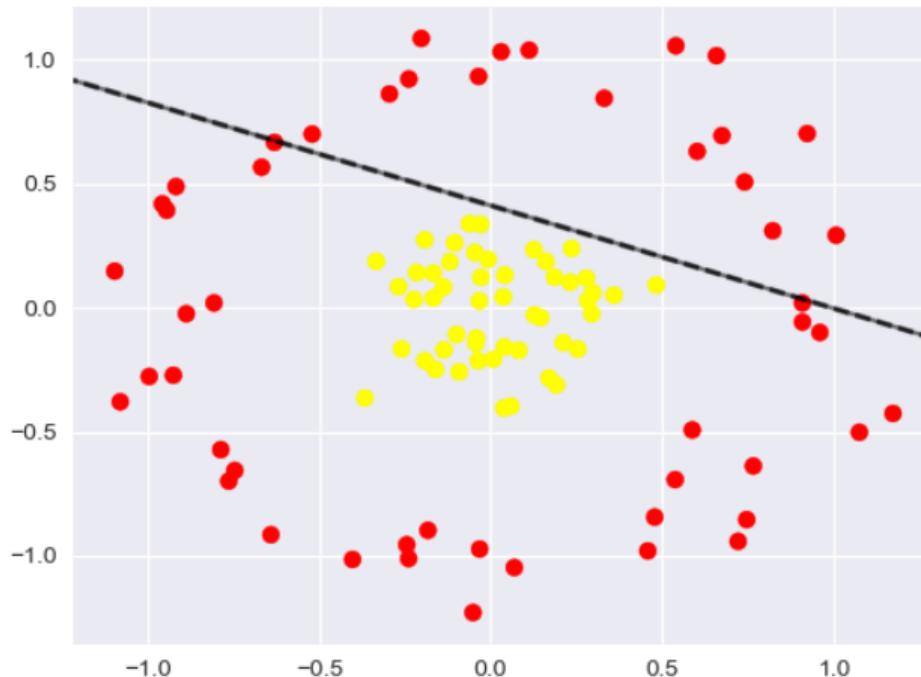


SVM - Max. margin shown in solid line





SVM - Max. margin shown in solid line



What to do here? Data is not linearly separable!

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

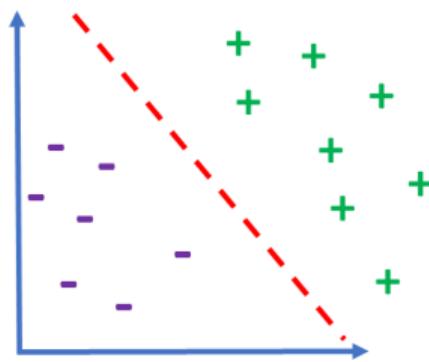
- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Problem Statement

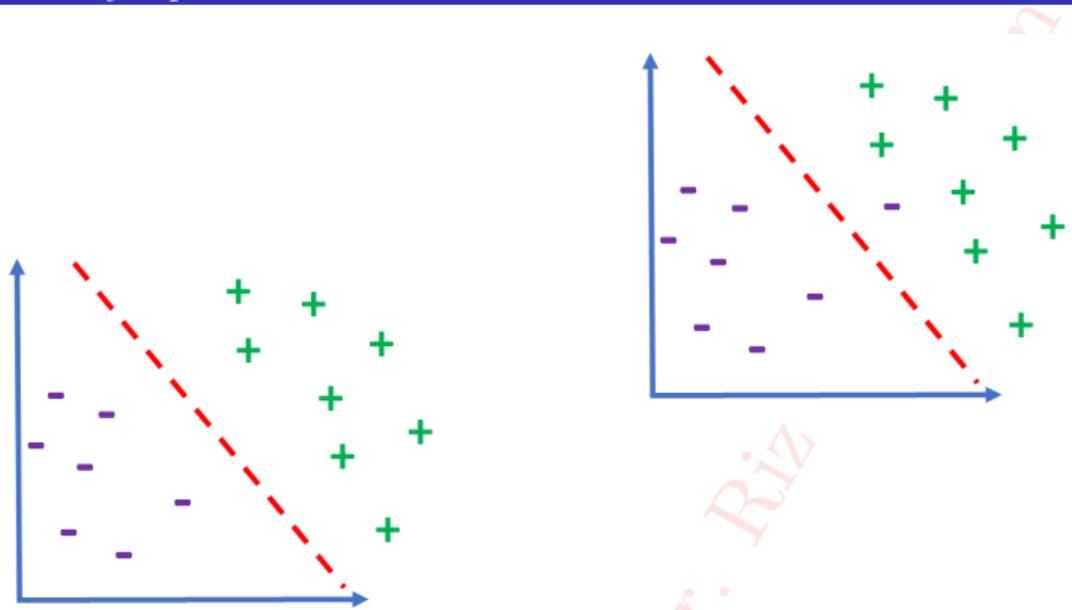
Non-linearly separable



(c)Dr. Rizwan A Khan

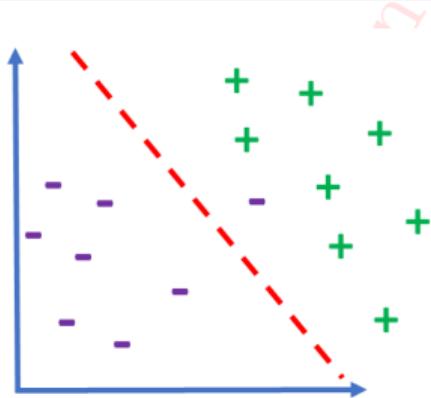
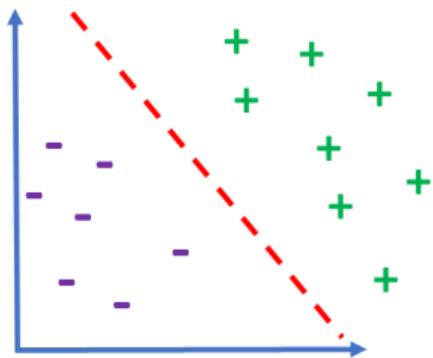
Problem Statement

Non-linearly separable



Problem Statement

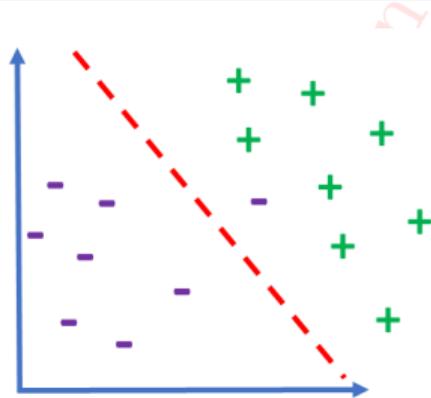
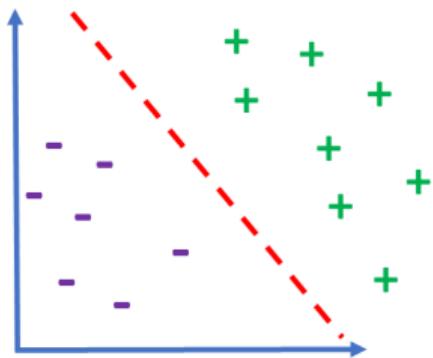
Non-linearly separable



- What to do, since SVM find linear classification boundary? SVMs are linear classifiers (a line in 2 dimensions, a plane in 3 dimensions, a $n - 1$ dimensional hyperplane in n dimensions).
- Some probabilities:

Problem Statement

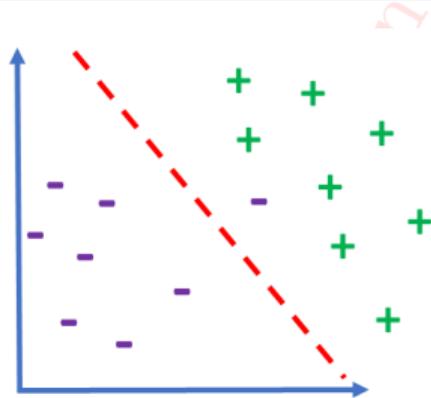
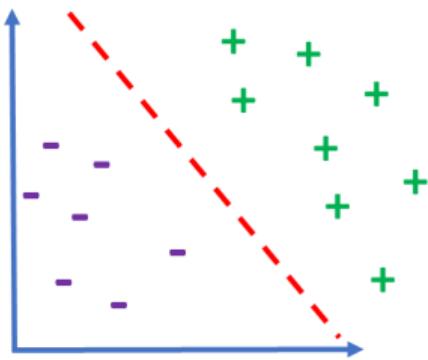
Non-linearly separable



- What to do, since SVM find linear classification boundary? SVMs are linear classifiers (a line in 2 dimensions, a plane in 3 dimensions, a $n - 1$ dimensional hyperplane in n dimensions).
- Some probabilities:
 - ➊ Probably its a outlier!

Problem Statement

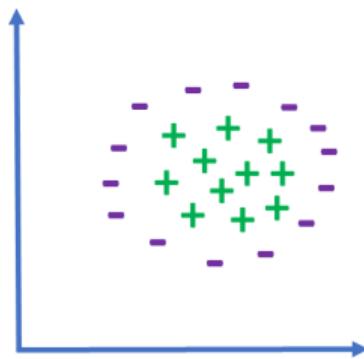
Non-linearly separable



- What to do, since SVM find linear classification boundary? SVMs are linear classifiers (a line in 2 dimensions, a plane in 3 dimensions, a $n - 1$ dimensional hyperplane in n dimensions).
- Some probabilities:
 - ➊ Probably its a outlier!
 - ➋ or find linear decision boundary while minimizing some cost function that penalizes for training error.

Problem Statement

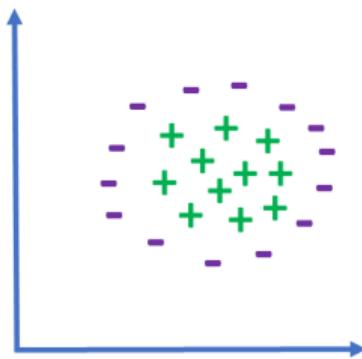
Non-linearly separable



(c)Dr. Rizwan A Khan

Problem Statement

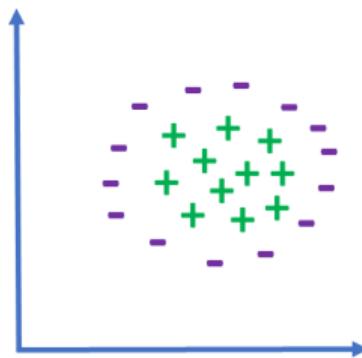
Non-linearly separable



- In previous example, one example seemed outlier and solution was proposed, but in this case it is impossible to come up with linear classifier.

Problem Statement

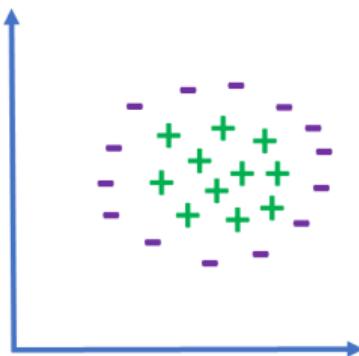
Non-linearly separable



- In previous example, one example seemed outlier and solution was proposed, but in this case it is impossible to come up with linear classifier.

- As it seems impossible to use SVM / linear classifier, should we just remove SVM from machine learning toolkit?

Non-linearly separable

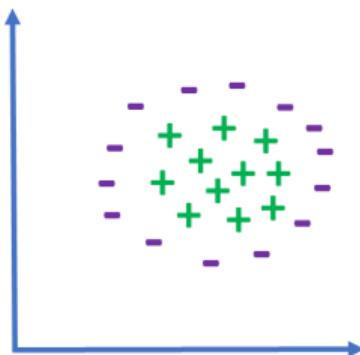


- In previous example, one example seemed outlier and solution was proposed, but in this case it is impossible to come up with linear classifier.

- As it seems impossible to use SVM / linear classifier, should we just remove SVM from machine learning toolkit?
- No, there is a little trick that can be done to transform data (change the data point without changing the data point) in a such away that it become linearly separable.

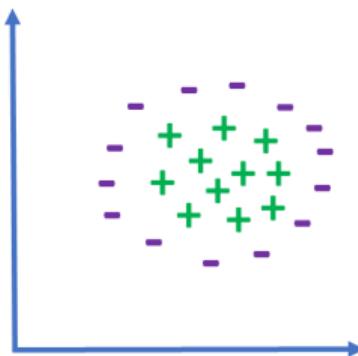
Problem Statement

Non-linearly separable



- In previous example, one example seemed outlier and solution was proposed, but in this case it is impossible to come up with linear classifier.

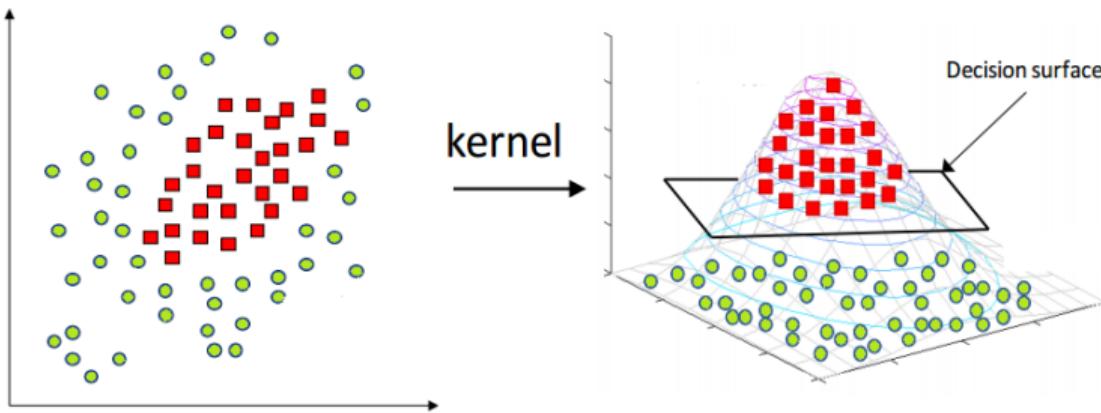
- As it seems impossible to use SVM / linear classifier, should we just remove SVM from machine learning toolkit?
- No, there is a little trick that can be done to transform data (change the data point without changing the data point) in a such away that it become linearly separable.
- Define function Φ that will take data point and change its dimension (in our example there are two dimension).



- In previous example, one example seemed outlier and solution was proposed, but in this case it is impossible to come up with linear classifier.

- As it seems impossible to use SVM / linear classifier, should we just remove SVM from machine learning toolkit?
- No, there is a little trick that can be done to transform data (change the data point without changing the data point) in a such away that it become linearly separable.
- Define function Φ that will take data point and change its dimension (in our example there are two dimension).
- For example, we can transform data from our example in three dimensions using:
$$\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle$$
where: x_1 and x_2 are dimension of same vector

Intuition

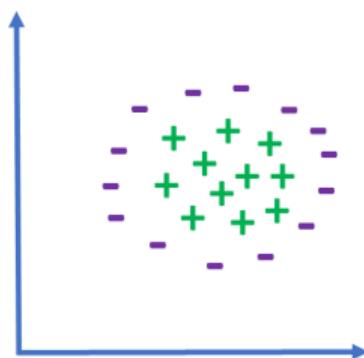


Video showing XOR data from 2D to 3D.³

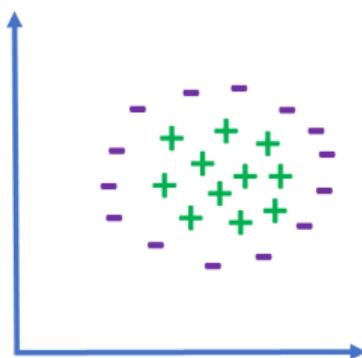
- There is a little trick that can be done to transform data (change the data point without changing the data point) in a such away that it become linearly separable.
- Define function Φ that will take data point and change its dimension.

³<https://youtu.be/5KIYu3zKvqo>

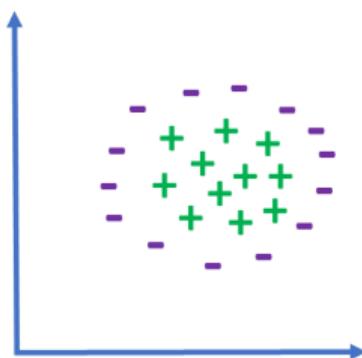
Non-linearly separable



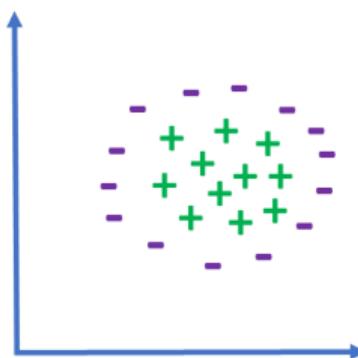
- $\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle$
- There isn't any new information!



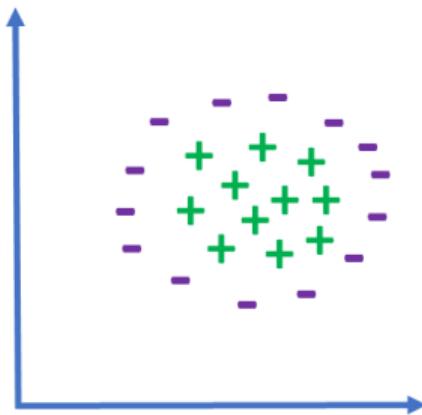
- $\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle$
- There isn't any new information!
- **Reminder:** Quadratic programming problem form:
$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$



- $\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle$
- There isn't any new information!
- **Reminder:** Quadratic programming problem form:
$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$
- **Reminder:** This optimization / finding extremum of a function depends only on dot products of pairs of samples, given by $\bar{X}_i^T \cdot \bar{X}_j$.

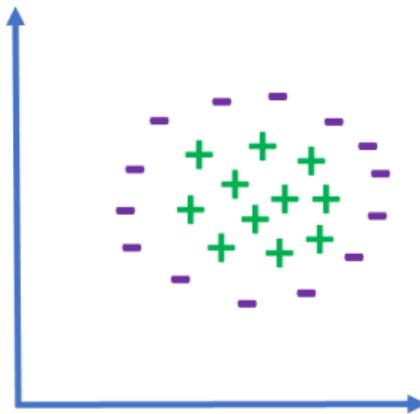


- $\Phi(x) = \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle$
- There isn't any new information!
- **Reminder:** Quadratic programming problem form:
$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$
- **Reminder:** This optimization / finding extremum of a function depends only on dot products of pairs of samples, given by $\bar{X}_i^T \cdot \bar{X}_j$.
- What will be $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given $\Phi(x)$



- **QUIZ:** What will be $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given $\Phi(x)$?

- Consider X_i as x and X_j as y , for the ease of notations:
 $\Phi(x)^T \Phi(y) =$



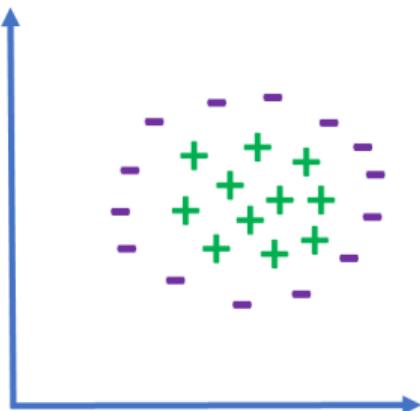
- QUIZ:** What will be
 $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given
 $\Phi(x)$?

Non-linearly separable

- Consider X_i as x and X_j as y , for the ease of notations:

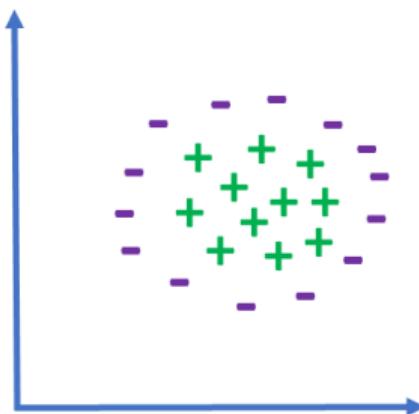
$$\Phi(x)^T \Phi(y) =$$

$$\begin{aligned} & \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle^T \cdot \langle y_1^2, y_2^2, \sqrt{2}y_1y_2 \rangle \\ & x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \quad (22) \\ & (x_1y_1 + x_2y_2)^2 \end{aligned}$$



- QUIZ:** What will be $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given $\Phi(x)$?

Non-linearly separable



- Consider X_i as x and X_j as y , for the ease of notations:

$$\Phi(x)^T \Phi(y) =$$

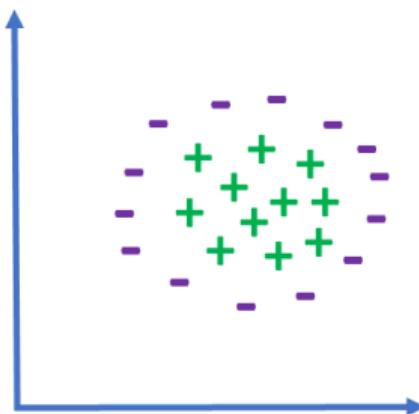
$$\begin{aligned} & \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle^T \cdot \langle y_1^2, y_2^2, \sqrt{2}y_1y_2 \rangle \\ & x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \quad (22) \\ & (x_1y_1 + x_2y_2)^2 \end{aligned}$$

- or this is equal to:

$$(x^T y)^2 \quad (23)$$

- QUIZ:** What will be $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given $\Phi(x)$?

Non-linearly separable



- Consider X_i as x and X_j as y , for the ease of notations:

$$\Phi(x)^T \Phi(y) =$$

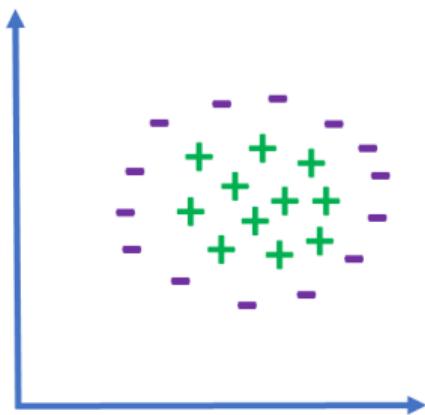
$$\begin{aligned} & \langle x_1^2, x_2^2, \sqrt{2}x_1x_2 \rangle^T \cdot \langle y_1^2, y_2^2, \sqrt{2}y_1y_2 \rangle \\ & x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \quad (22) \\ & (x_1y_1 + x_2y_2)^2 \end{aligned}$$

- or this is equal to:

$$(x^T y)^2 \quad (23)$$

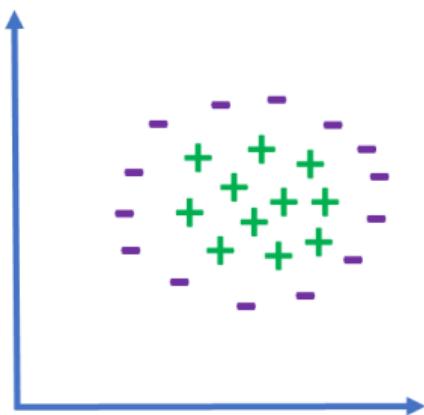
- QUIZ:** What will be $\Phi(\bar{X}_i)^T \cdot \Phi(\bar{X}_j)$ for a given $\Phi(x)$?

- So, dot product becomes square of last dot product.

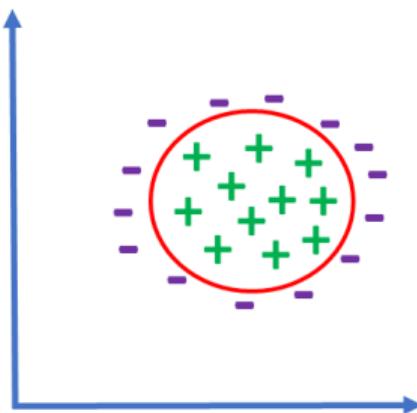


- Refer [Equation 23](#), its particular form of equation of circle in matrix form.

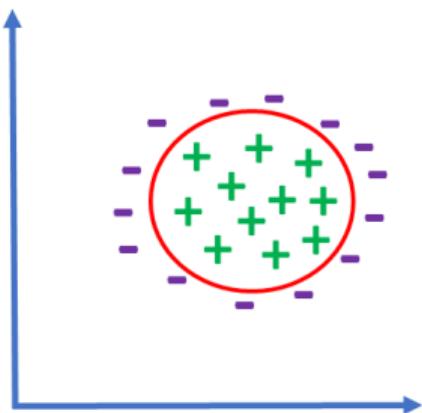
(c)Dr. Rizwan Khan



- Refer [Equation 23](#), its particular form of equation of circle in matrix form.
- Somehow, the notion of similarity ($\bar{X}_i^T \cdot \bar{X}_j$) is transformed to notion of circle where some points remain inside the circle while others don't.



- Refer [Equation 23](#), its particular form of equation of circle in matrix form.
- Somehow, the notion of similarity ($\bar{X}_i^T \cdot \bar{X}_j$) is transformed to notion of circle where some points remain inside the circle while others don't.
- So, data got transformed from 2D to 3D (without any additional information) in such a way that now it can be separated by hyperplane.



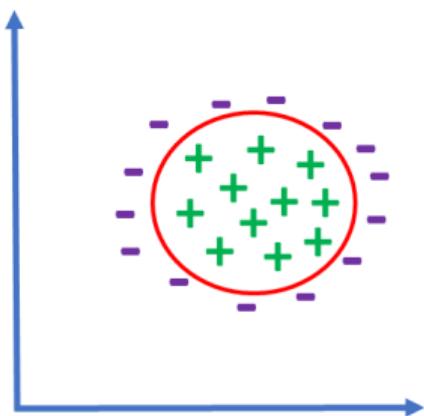
- Refer [Equation 23](#), its particular form of equation of circle in matrix form.
- Somehow, the notion of similarity ($\bar{X}_i^T \cdot \bar{X}_j$) is transformed to notion of circle where some points remain inside the circle while others don't.
- So, data got transformed from 2D to 3D (without any additional information) in such a way that now it can be separated by hyperplane.
- This little trick of projecting data into higher dimension space to make it separable is called [Kernel trick](#).

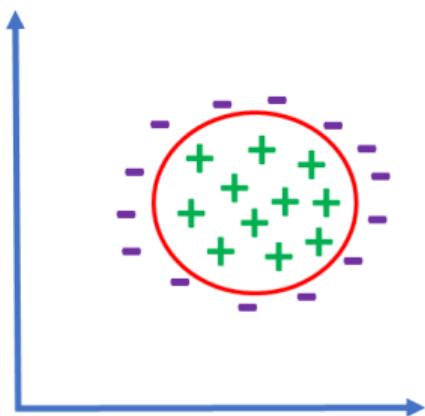
- Coming back to this equation:
Quadratic programming problem form:

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$

and Equation 23:

$$(x^T y)^2$$





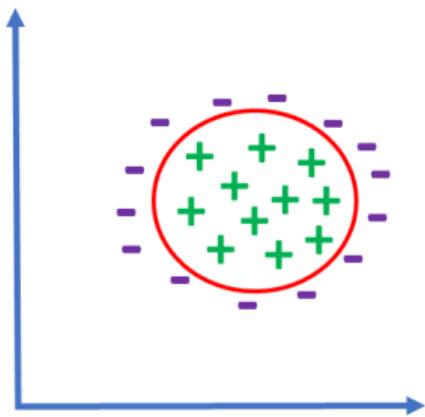
- Coming back to this equation:
Quadratic programming problem form:

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$

and Equation 23:

$$(x^T y)^2$$

- This signifies that data points don't need to be transformed separately, but rather its just a dot product squared! **So we actually never used Φ .**



- Coming back to this equation:
Quadratic programming problem form:

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$

and Equation 23:

$$(x^T y)^2$$

- This signifies that data points don't need to be transformed separately, but rather its just a dot product squared! **So we actually never used Φ .**
- That's the beauty of mathematics and SVM.**

Kernel Trick

- $$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{X}_i \cdot \bar{X}_j$$

- This equation can now be written as:

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\bar{X}_i \cdot \bar{X}_j) \quad (24)$$

Where K is Kernel function, that takes X_i and X_j and returns scalar, the inner product (generalization of the dot product) between two points in a suitable feature space.

- Kernel functions allow to inject domain knowledge into classifier.

Kernel Trick

Kernel Trick

- Gaussian Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}} \quad (25)$$

- Polynomial

$$K(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i)^T \mathbf{x}_j + c)^p \quad (26)$$

- Neural-net inspired!

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\kappa \mathbf{x}_i \mathbf{x}_j - \delta)^p \quad (27)$$

- Radial Basis

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{(\gamma ||\mathbf{x}_i - \mathbf{x}_j||^2)} \quad (28)$$

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Python for SVM

```
1 from __future__ import division, print_function
2 import numpy as np
3 from sklearn import datasets, svm
4 #from sklearn.cross_validation import train_test_split
5 from sklearn.model_selection import train_test_split
6 import matplotlib.pyplot as plt
7
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.ensemble import RandomForestClassifier, BaggingClassifier,
    AdaBoostClassifier, VotingClassifier
10
11 iris = datasets.load_iris()
12 X = iris.data[:,2:] # First two features, can take last two using [:,2:]
13 y = iris.target
14
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
    random_state=42)
```

Python for SVM

```
1 def evaluate_on_test_data(model=None):
2     predictions = model.predict(X_test)
3     correct_classifications = 0
4     for i in range(len(y_test)):
5         if predictions[i] == y_test[i]:
6             correct_classifications += 1
7     accuracy = 100*correct_classifications/len(y_test) #Accuracy as a
8     percentage
9     return accuracy
10
11 kernels = ('linear', 'poly', 'rbf')
12 accuracies = []
13 for index, kernel in enumerate(kernels):
14     model = svm.SVC(kernel=kernel)
15     model.fit(X_train, y_train)
16     acc = evaluate_on_test_data(model)
17     accuracies.append(acc)
18     print("{} % Test accuracy obtained with kernel = {}".format(acc, kernel))
```

Python for SVM

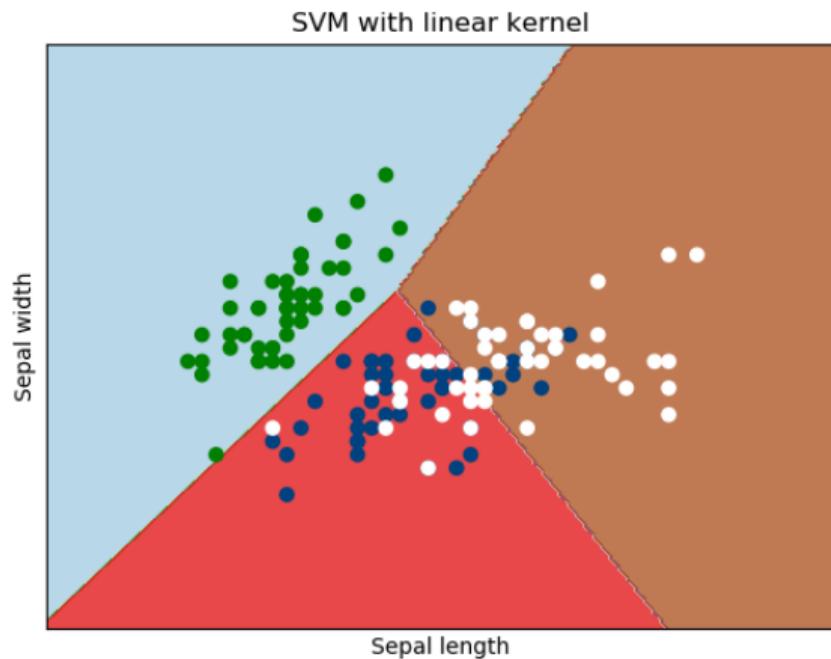
an

```
1 #Train SVMs with different kernels
2 svc = svm.SVC(kernel='linear').fit(X_train, y_train)
3 rbf_svc = svm.SVC(kernel='rbf', gamma=0.7).fit(X_train, y_train)
4 poly_svc = svm.SVC(kernel='poly', degree=9).fit(X_train, y_train)
5
6
7
8 #Create a mesh to plot in
9 h = .02 # step size in the mesh
10 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
11 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
12 xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
13                      np.arange(y_min, y_max, h))
14
15 #Define title for the plots
16 titles = ['SVM with linear kernel',
17            'SVM with RBF kernel',
18            'SVM with polynomial (degree 9) kernel']
```

Python for SVM

```
1 for i, clf in enumerate((svc, rbf_svc, poly_svc)):
2     # Plot the decision boundary. For that, we will assign a color to each
3     # point in the mesh [x_min, x_max]x[y_min, y_max].
4     plt.figure(i)
5
6     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
7     # Put the result into a color plot
8     Z = Z.reshape(xx.shape)
9     plt.contourf(xx, yy, Z, cmap=plt.cm.Paired, alpha=0.8)
10
11    # Plot also the training points
12    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.ocean)
13    plt.xlabel('Sepal length')
14    plt.ylabel('Sepal width')
15    plt.xlim(xx.min(), xx.max())
16    plt.ylim(yy.min(), yy.max())
17    plt.xticks(())
18    plt.yticks(())
19    plt.title(titles[i])
20 plt.show()
```

SVM Visualization



SVM with Linear Kernel (No transformation)

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j + c$$

Introduction
ooooooo

Decision Rule
oooooo

Constraints
oooooooooooo

Optimal Margin Classifier
oooooooooooo

Kernel Trick
oooooooooooo

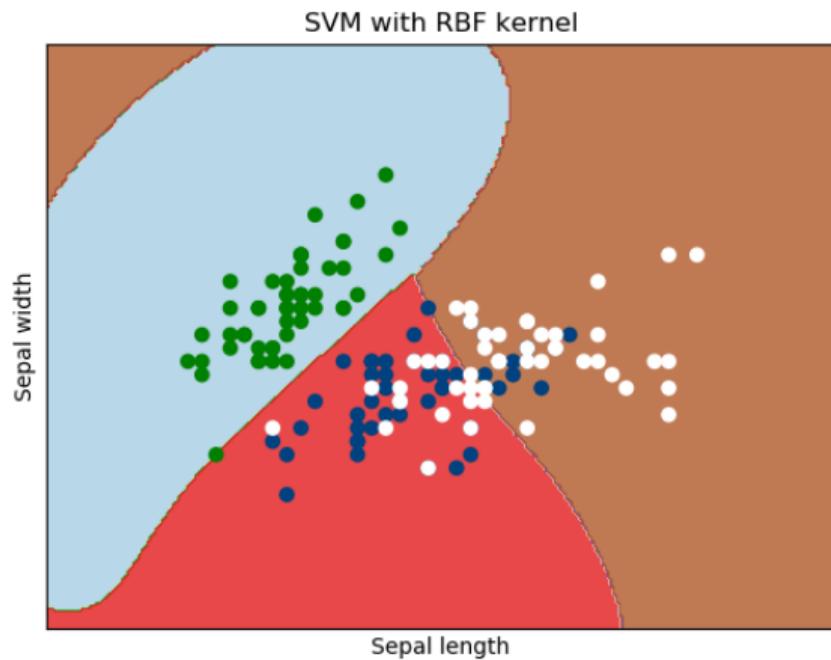
Python
ooooooo●oooo

Soft Margin SVM
ooooooo

Lagrange Optimization
oooooooooooo

Results

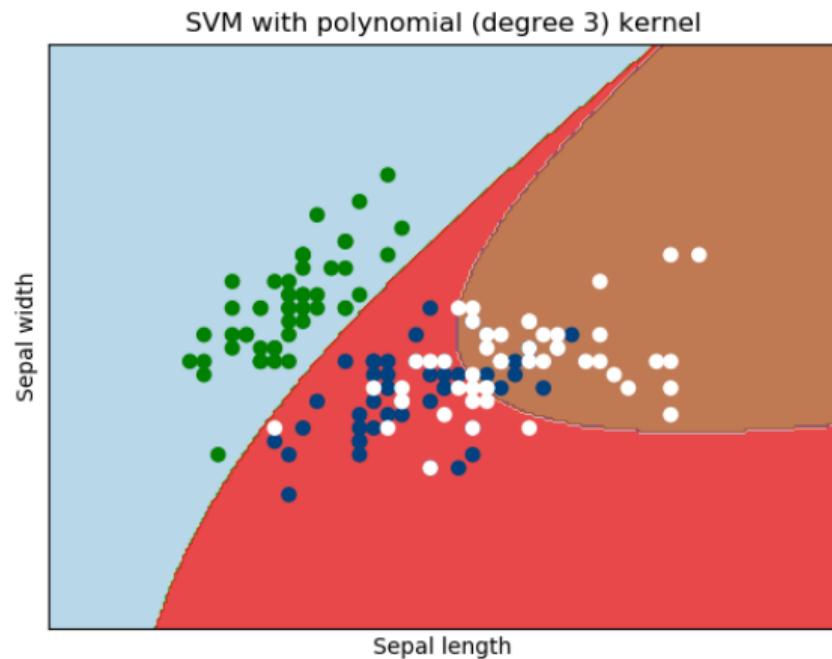
SVM Visualization



SVM with RBF

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$$

SVM Visualization

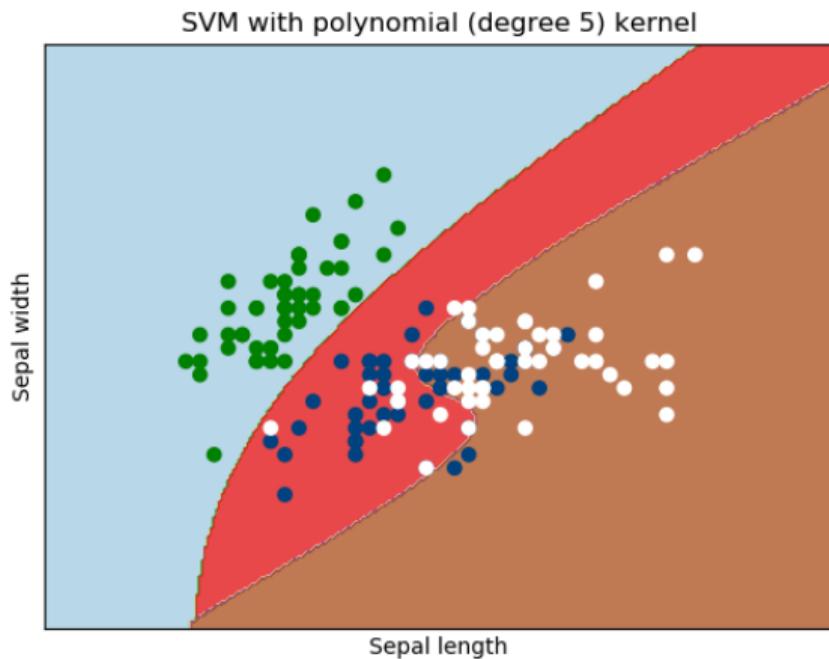


Data

SVM with Polynomial Kernel (Degree 3)

$$K(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i)^T \mathbf{x}_j + c)^p$$

SVM Visualization



Data

SVM with Polynomial Kernel (Degree 5)

$$K(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i)^T \mathbf{x}_j + c)^p$$

Introduction
ooooooo

Decision Rule
ooooo

Constraints
oooooooooooo

Optimal Margin Classifier
oooooooooooo

Kernel Trick
oooooooooooo

Python
oooooo•oooo

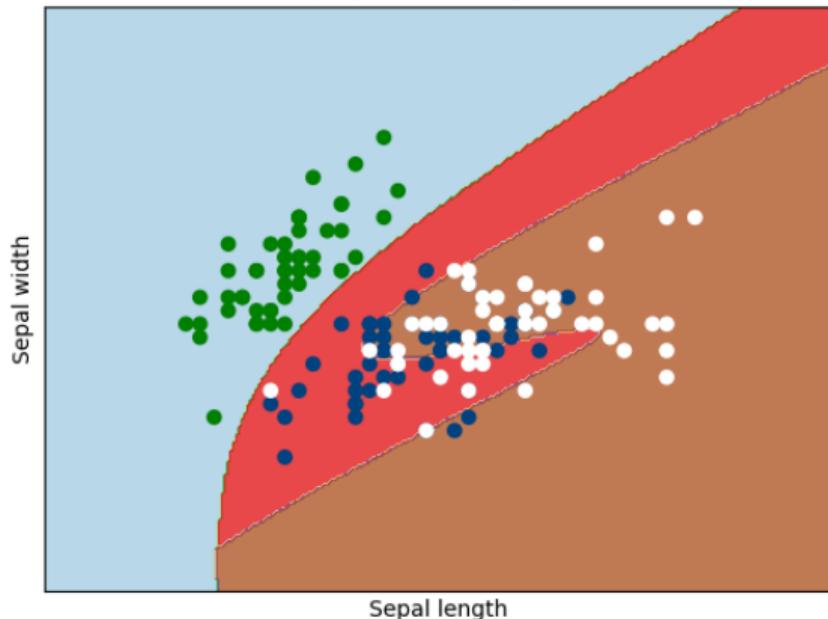
Soft Margin SVM
ooooooo

Lagrange Optimization
oooooooooooo

Results

SVM Visualization

SVM with polynomial (degree 7) kernel

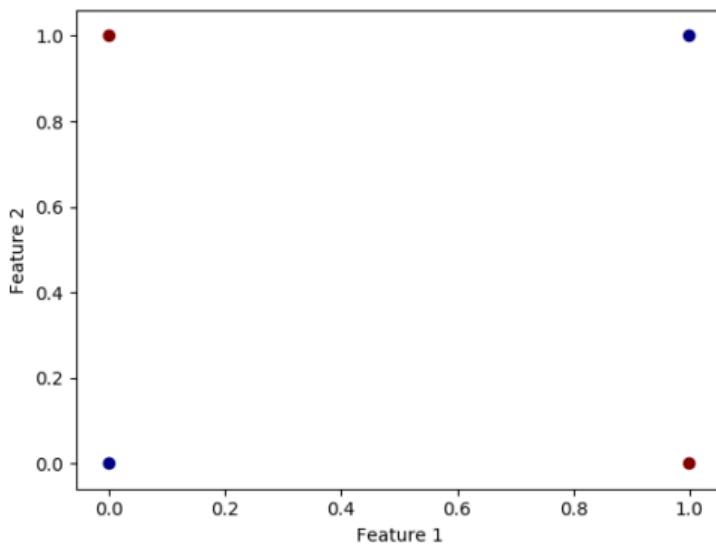


SVM with Polynomial Kernel (Degree 7)

$$K(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i)^T \mathbf{x}_j + c)^p$$

XOR problem visualization

XOR problem



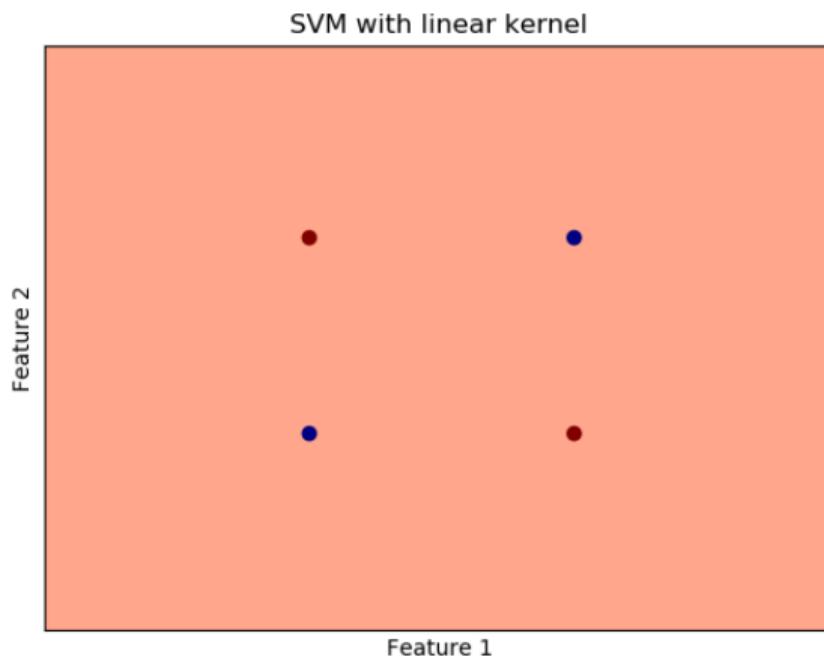
Video ⁴

- XOR data.
- Problem is **non-linearly separable** in the given feature space.

⁴<https://youtu.be/5KIYu3zKvqo>

XOR problem visualization

XOR problem

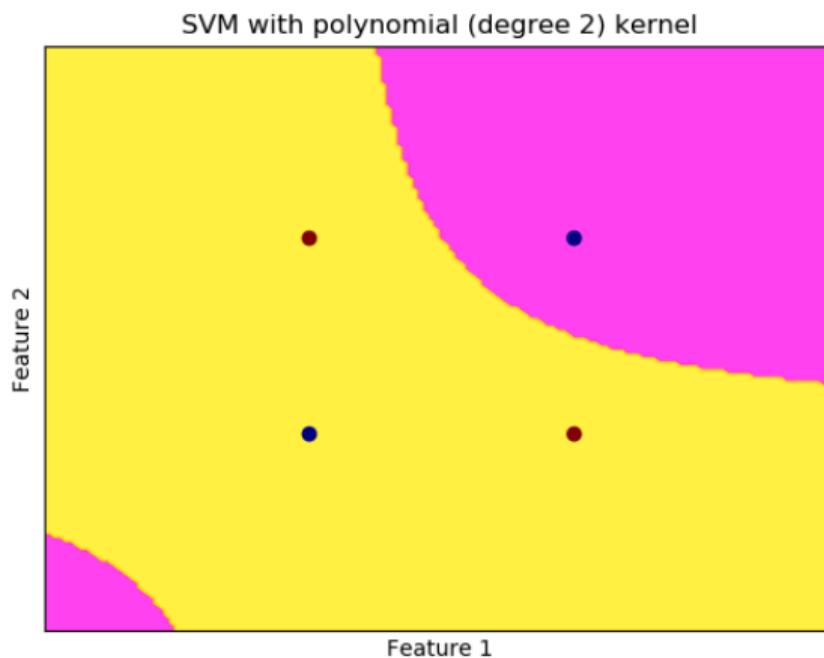


ad

SVM with Linear Kernel (No transformation)

XOR problem visualization

XOR problem

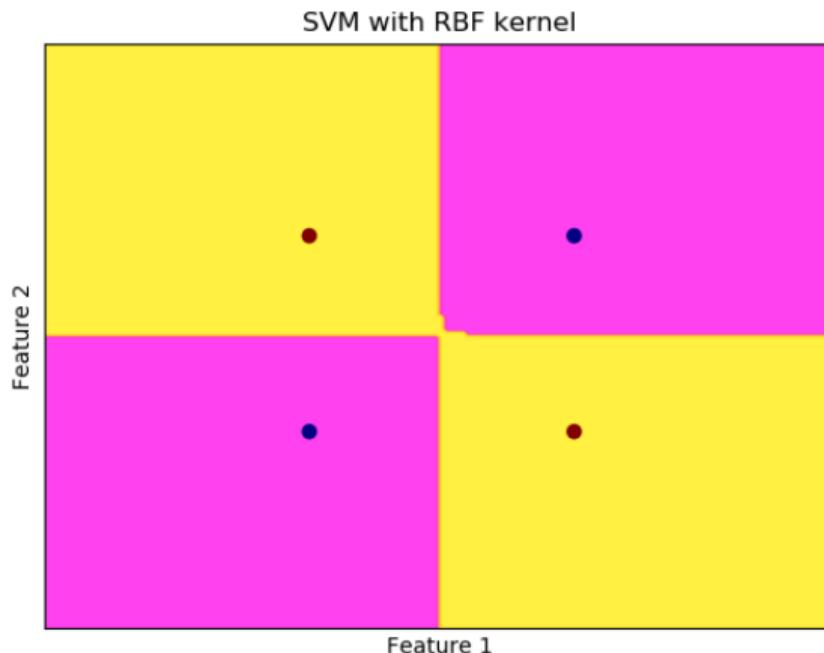


an

SVM with Polynomial Kernel

XOR problem visualization

XOR problem



drkhan

SVM with RBF Kernel

Conclusion

Conclusion

- Powerful theoretical grounds.
- Global, unique solution (convex optimization function).
- Performance depends on choice of kernel and parameters.
- Training is memory-intensive.
- Complexity dependent on the number of support vectors.

Video⁵

⁵https://youtu.be/FxLIsbnp_5c

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

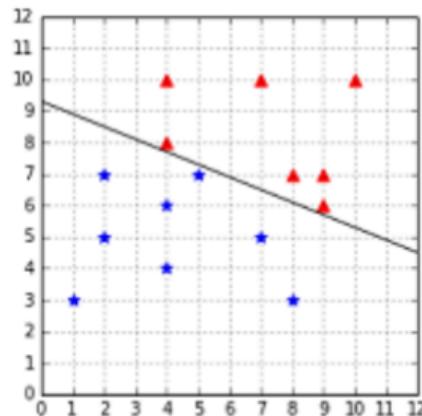
8 Lagrange Optimization

- Function Visualization
- Function Optimization

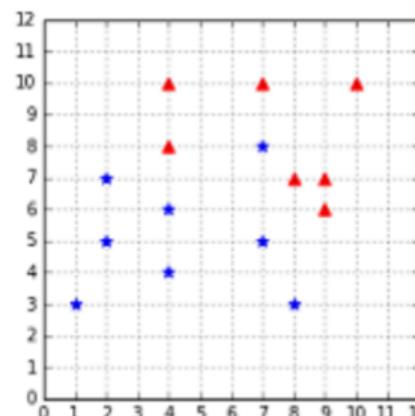
Need for soft Margin SVM

Need for soft Margin SVM

- Real-life data is often noisy, thus linear separability is an issue.
 - Even when the data is linearly separable, the outlier can be closer to the other examples than most of the examples of its class, thus reducing the margin, or it can be among the other examples and break linear separability.



Outlier reducing the margin



Outlier breaks linear separability

- In this case, there is no solution to the optimization problem solved earlier.

⁶Image taken from SVM Succinctly

Soft margin
Slack variable

- In 1995, Vapnik and Cortes introduced a modified version of the original SVM that allows the classifier to make some mistakes.
- The goal is now not to make zero classification mistakes, but to **make as few mistakes as possible**.
- Constraints of the optimization problem was modified, so the constraint (refer Eq 4)

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1$$

becomes

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1 - \xi_i \quad (29)$$

where: ξ = slack variable and $\forall i \geq 0$.

- In 1995, Vapnik and Cortes introduced a modified version of the original SVM that allows the classifier to make some mistakes.
- The goal is now not to make zero classification mistakes, but to **make as few mistakes as possible**.
- Constraints of the optimization problem was modified, so the constraint (refer Eq 4)

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1$$

becomes

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1 - \xi_i \quad (29)$$

where: ξ = slack variable and $\forall i \geq 0$.

- ξ is subtracted from 1, in order to make it possible to satisfy constraint.

- The problem is that we could choose a huge value of ξ for every example, and all the constraints will be satisfied.
- To avoid this, we need to modify the objective function (refer Eq. 14 for objective function of hard margin SVM) to penalize the choice of a big ξ :

$$\operatorname{argmin}_{W,b,\xi} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i \quad (30)$$

subject to

$$y_i(\bar{W} \cdot \bar{X}_i + b) \geq 1 - \xi_i$$

and $\xi_i \geq 0 \forall i$

Soft margin

Slack variable

- Consider:

$$\operatorname{argmin}_{W,b,\xi} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i$$

The slack variable ξ_i allows the input xi to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such “slack”.

- How to select hyper-parameter C ?

Soft margin

Slack variable

- Consider:

$$\operatorname{argmin}_{W,b,\xi} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i$$

The slack variable ξ_i allows the input xi to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such “slack”.

- How to select hyper-parameter C ?

Value of C (hyper-parameter tuning)

- If C is very large (penalty is higher), the SVM becomes very strict and tries to classify all data points correctly.

Soft margin

Slack variable

- Consider:

$$\operatorname{argmin}_{W,b,\xi} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i$$

The slack variable ξ_i allows the input xi to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such “slack”.

- How to select hyper-parameter C ?

Value of C (hyper-parameter tuning)

- If C is very large (penalty is higher), the SVM becomes very strict and tries to classify all data points correctly.
- If C is very small, the SVM becomes very loose and may “sacrifice” some points to obtain a simpler solution.

- Consider:

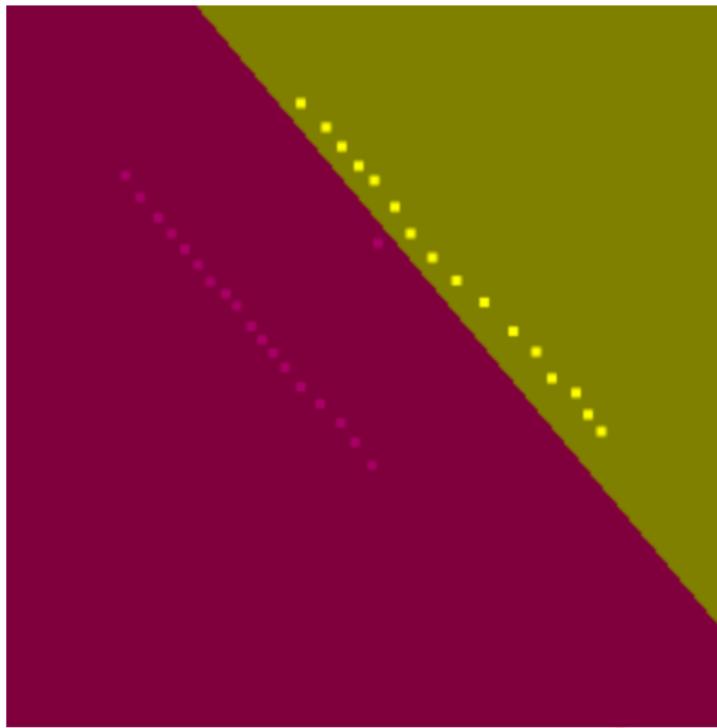
$$\operatorname{argmin}_{W,b,\xi} \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i$$

The slack variable ξ_i allows the input xi to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such “slack”.

- How to select hyper-parameter C ?

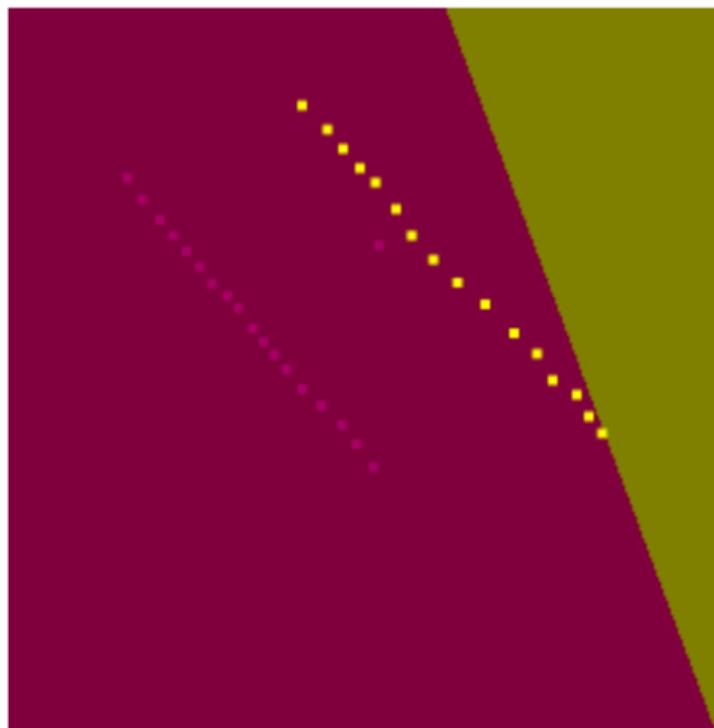
Value of C (hyper-parameter tuning)

- If C is very large (penalty is higher), the SVM becomes very strict and tries to classify all data points correctly.
- If C is very small, the SVM becomes very loose and may “sacrifice” some points to obtain a simpler solution.
- Usually telescopic / grid search is applied to find best C for the given dataset.



This image corresponds to large value of C .

⁷Demo images from LIBSVM website: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>



This image corresponds to small value of C , over-simplification of solution.

⁷Demo images from LIBSVM website: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>



*7

This image corresponds to appropriate value of C given dataset (maximizing margin, sacrificing few (one data point) to obtain wider margin / better generalization).

⁷Demo images from LIBSVM website: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

Section Contents

1 Introduction

- Reference Books
- Intuition - Decision Boundary

2 Decision Rule

- SVM - Decision Boundary

3 Constraints

- Distance b/w boundary lines
- Summary

4 Optimal Margin Classifier

- Summary
- Ocular proof

5 Kernel Trick

- Problem Statement

• Kernel Trick

6 Python

- Code
- Results
- XOR problem visualization
- Conclusion

7 Soft Margin SVM

- Need for soft Margin SVM
- Soft margin
- Ocular Proof

8 Lagrange Optimization

- Function Visualization
- Function Optimization

Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.

(c)Dr. Rizwan A Khan

Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)

Function Visualization

Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.

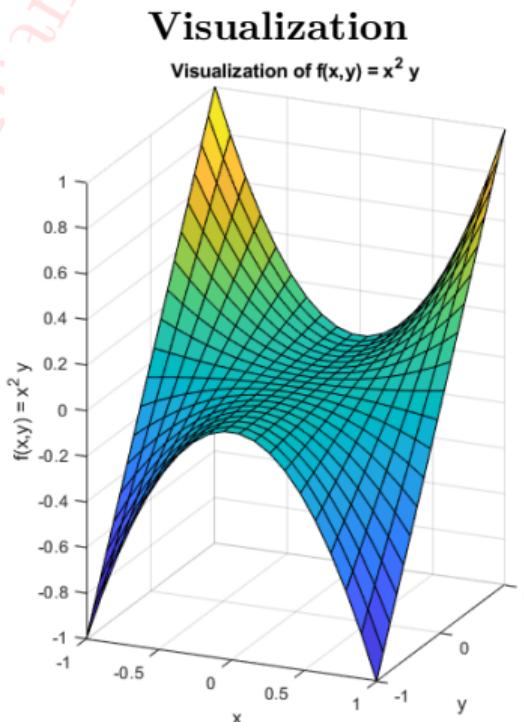
- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)

Visualization

(c)Dr. Rizwan A Khan

Optimizing function with constraint

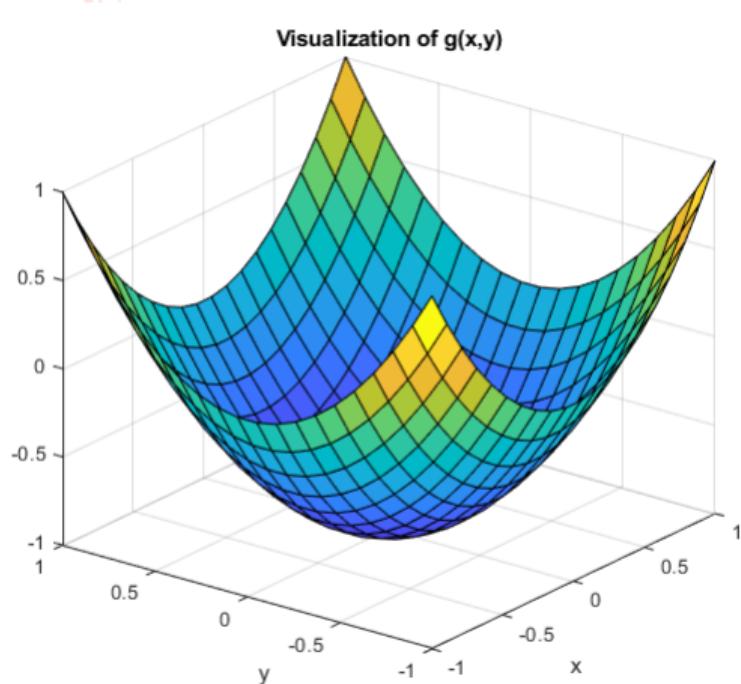
- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)



The function is in 3-D. To analytically examine the problem , we can use contour plot.

Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)



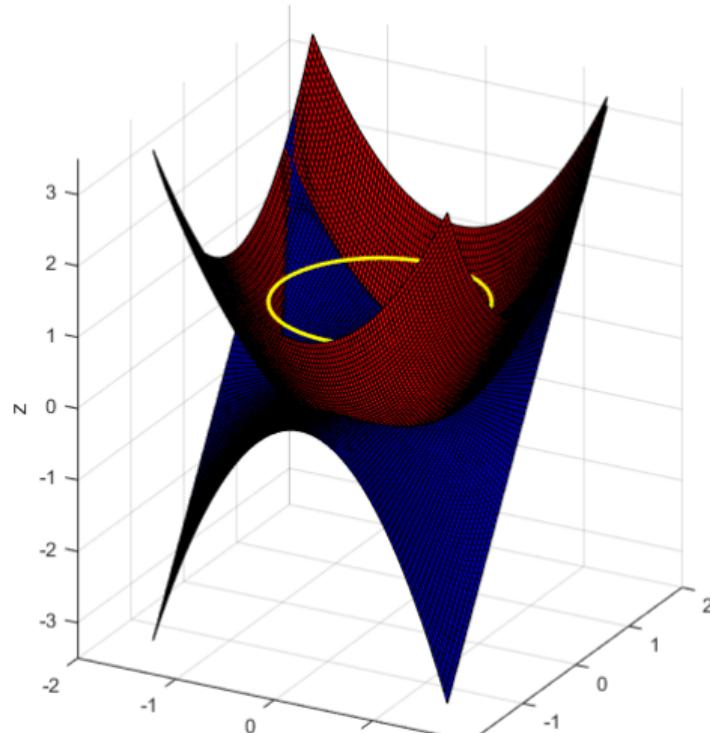
Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)

the set
action
that

Visualization

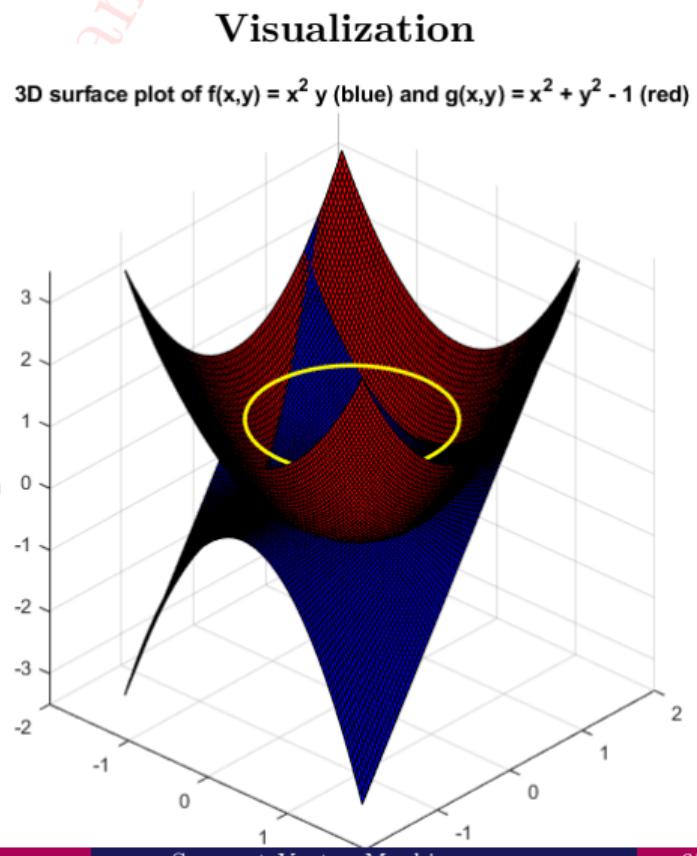
3D surface plot of $f(x,y) = x^2 y$ (blue) and $g(x,y) = x^2 + y^2 - 1$ (red)



Function Visualization

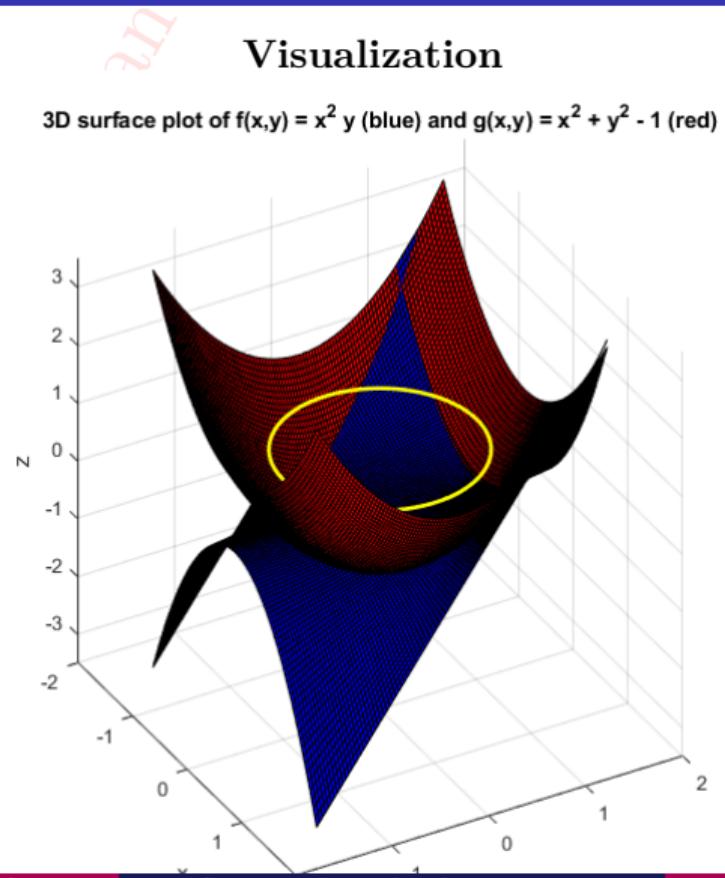
Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)



Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)



Function Visualization

Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.

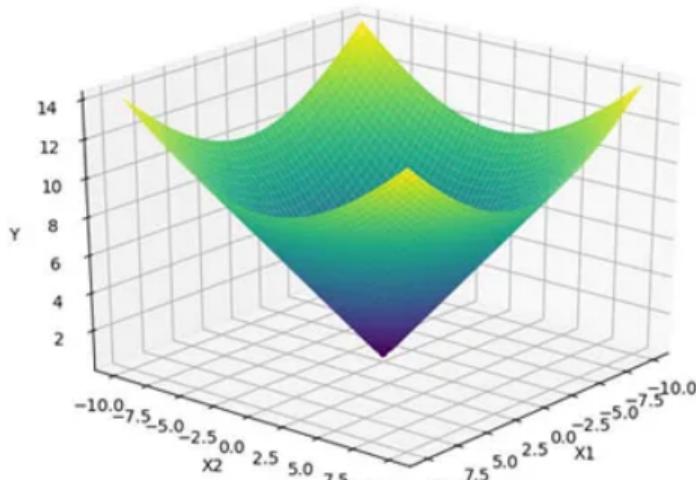
- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint $g(x, y)$ that $x^2 + y^2 = 1$ (unit circle)

Visualization

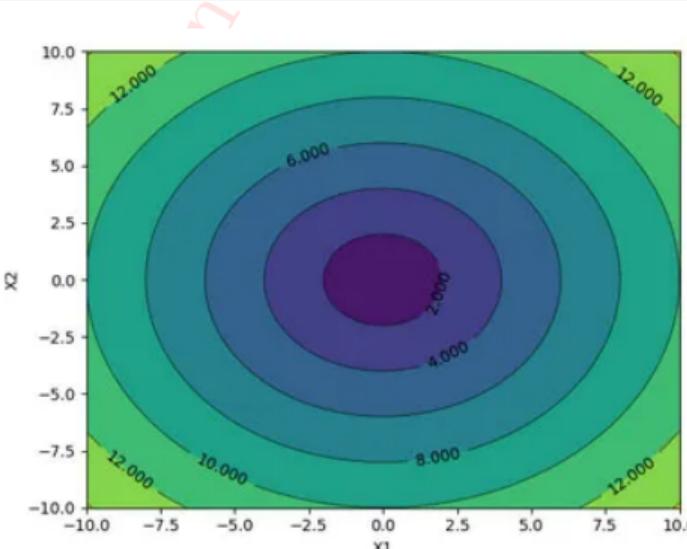
Contour plot

A contour plot is a graphical technique for representing a 3-D surface by plotting constant z slices, called contours, on a 2-D format. That is, lines are drawn for all possible pairs of (x, y) that produce same/constant output (z).

Contour Plot



3D Plot

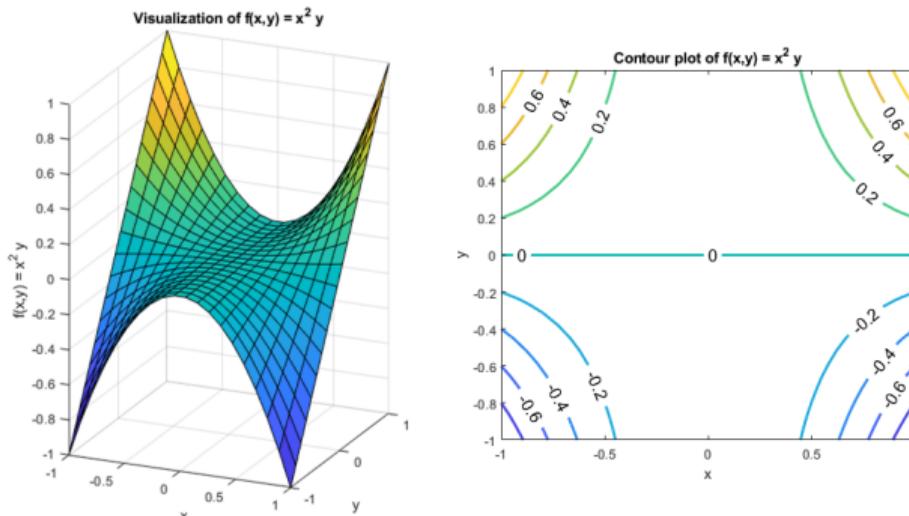


Contour Plot

Contour plot

A contour plot is a graphical technique for representing a 3-D surface by plotting constant z slices, called contours, on a 2-D format. That is, lines are drawn for all possible pairs of (x, y) that produce same/constant output (z).

Contour Plot



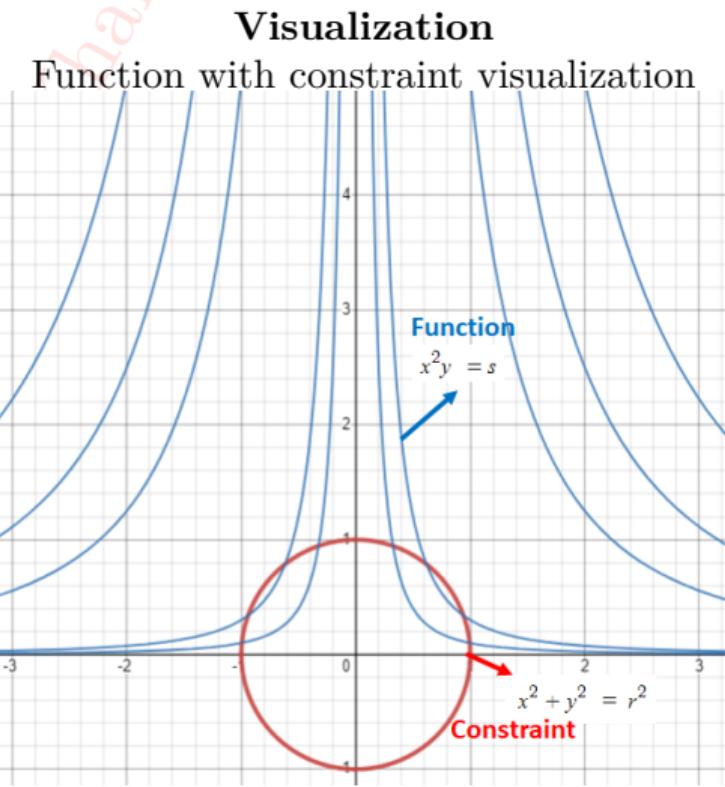
3D plot of function, along with corresponding contour plot of the problem in hand.

Contour plot

A contour plot is a graphical technique for representing a 3-D surface by plotting constant z slices, called contours, on a 2-D format. That is, lines are drawn for all possible pairs of (x, y) that produce same/constant output (z).

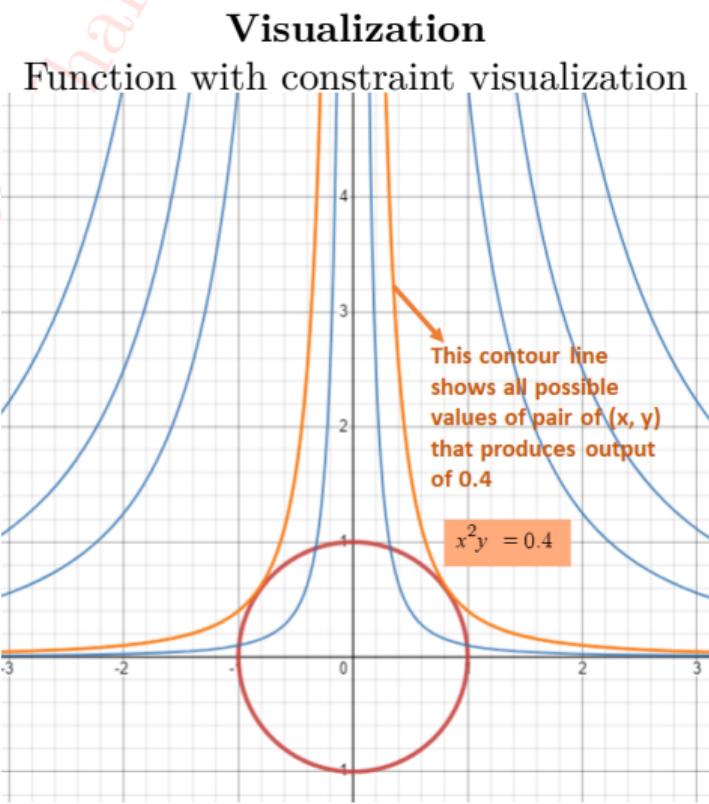
Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint that $x^2 + y^2 = 1$ (unit circle)

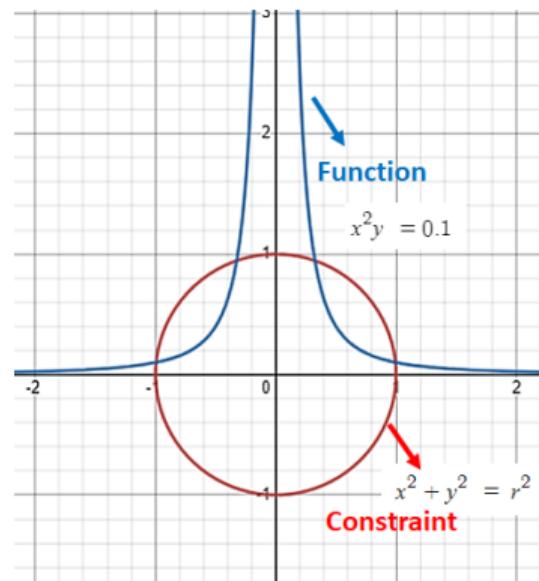


Optimizing function with constraint

- Maximize $f(x, y) = x^2y$ on the set $x^2 + y^2 = 1$.
 - Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint that $x^2 + y^2 = 1$ (unit circle)

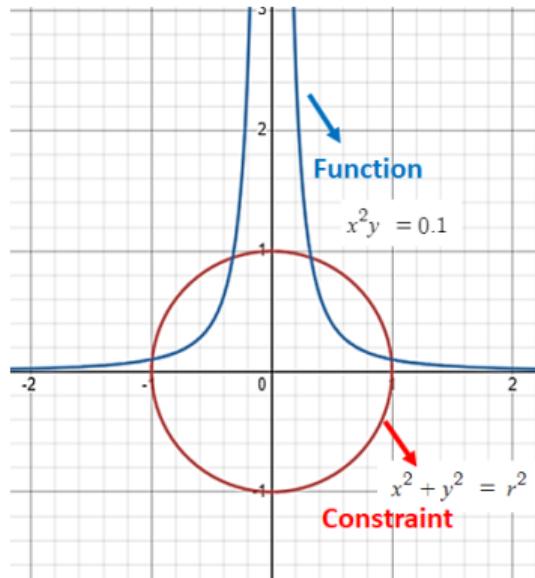


Optimizing function with constraint

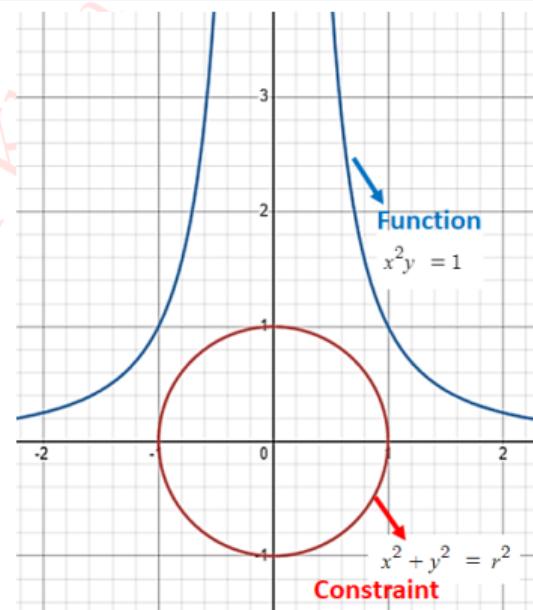


Function intersects with the constraint.
This means that this pair of (x, y)
satisfies constraint (four possible pair of
 (x, y) values), but visually we can
observe that they are maximum values.

Optimizing function with constraint



Function intersects with the constraint. This means that this pair of (x, y) satisfies constraint (four possible pair of (x, y) values), but visually we can observe that they are maximum values.



Function never intersects with the constraint.
 This means that this pair of (x, y) is off the constraint. It also shows that, as we are max. x^2 subject to constraint, we can never go as high as these values of (x, y) .

Optimizing function with constraint

- Here we are trying to

- Optimize multi-variable function

$$f(x, y) = x^2y$$

- with the constraint that $x^2 + y^2 = 1$
(unit circle)

Function Visualization

Optimizing function with constraint

- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint that $x^2 + y^2 = 1$
(unit circle)

Objective

To maximize function $f(x, y) = x^2y$ while satisfying constraint $x^2 + y^2 = 1$, is to find maximum value of pair of (x, y) or value of constant s to the point that afterwards its off the constraint.

This will only happen when the two functions ($f(x, y)$ & $g(x, y)$) are **tangent**.

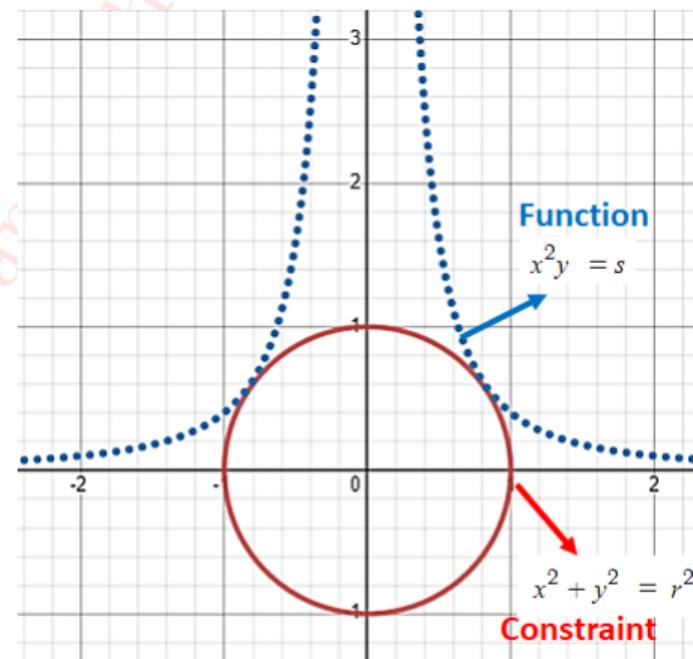
Optimizing function with constraint

- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint that $x^2 + y^2 = 1$ (unit circle)

Objective

To maximize function $f(x, y) = x^2y$ while satisfying constraint $x^2 + y^2 = 1$, is to find maximum value of pair of (x, y) or value of constant s to the point that afterwards its off the constraint.

This will only happen when the two functions ($f(x, y)$ & $g(x, y)$) are **tangent**.



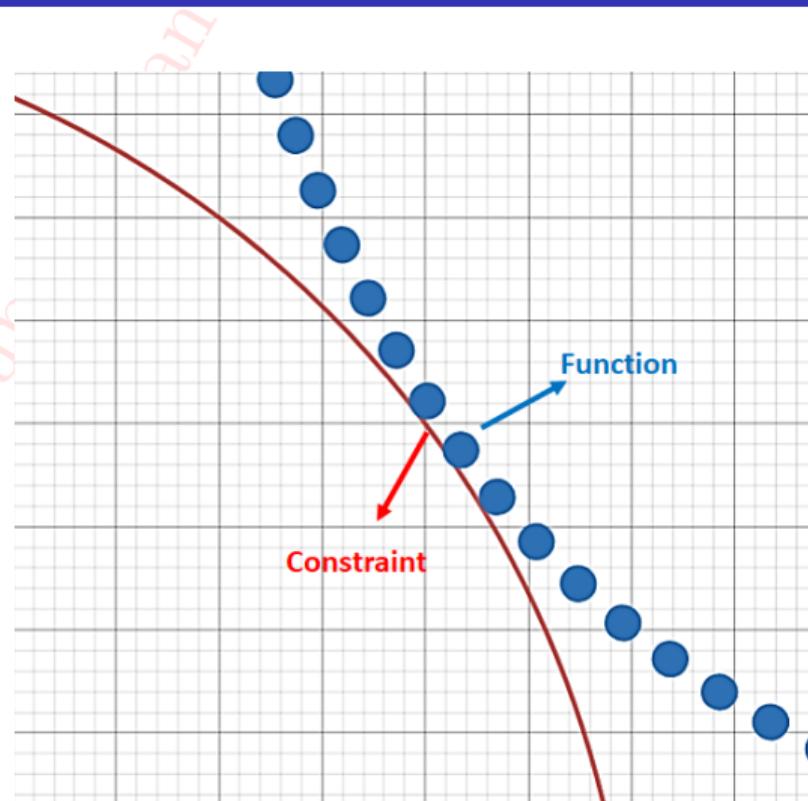
Optimizing function with constraint

- Here we are trying to
 - Optimize multi-variable function $f(x, y) = x^2y$
 - with the constraint that $x^2 + y^2 = 1$ (unit circle)

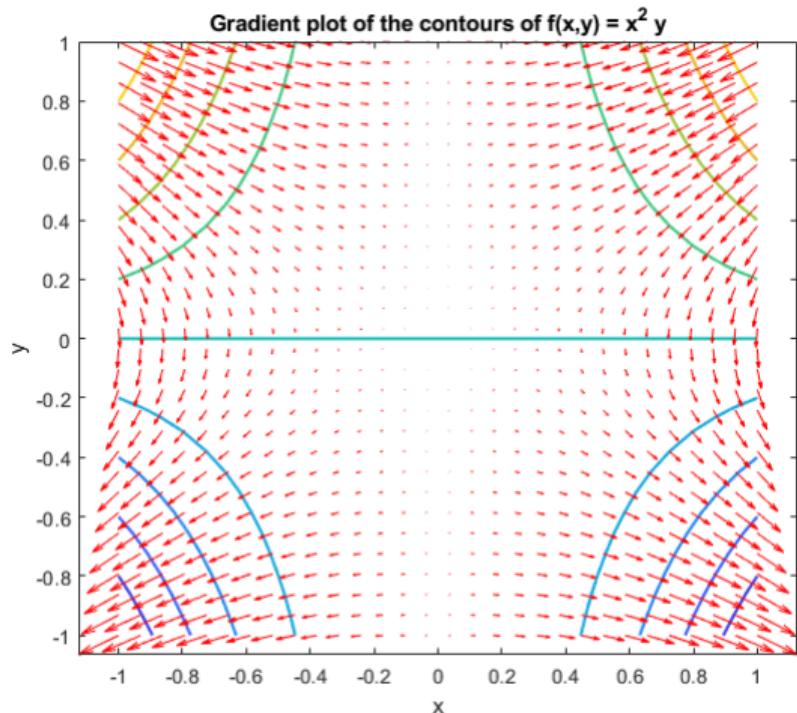
Objective

To maximize function $f(x, y) = x^2y$ while satisfying constraint $x^2 + y^2 = 1$, is to find maximum value of pair of (x, y) or value of constant s to the point that afterwards its off the constraint.

This will only happen when the two functions ($f(x, y)$ & $g(x, y)$) are **tangent**.



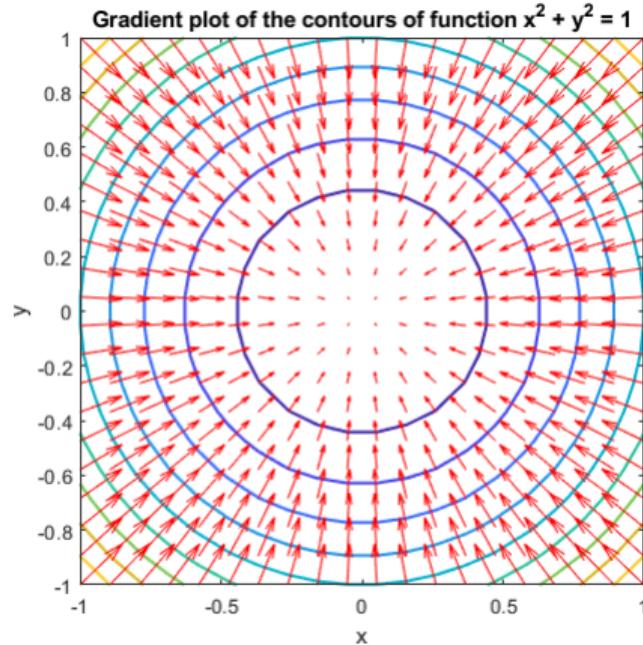
Optimizing function with constraint



The gradient of f or g evaluated at a point (x_0, y_0) always gives a vector perpendicular to the contour line passing through that point (as there is no change in value along contour line).

Function Optimization

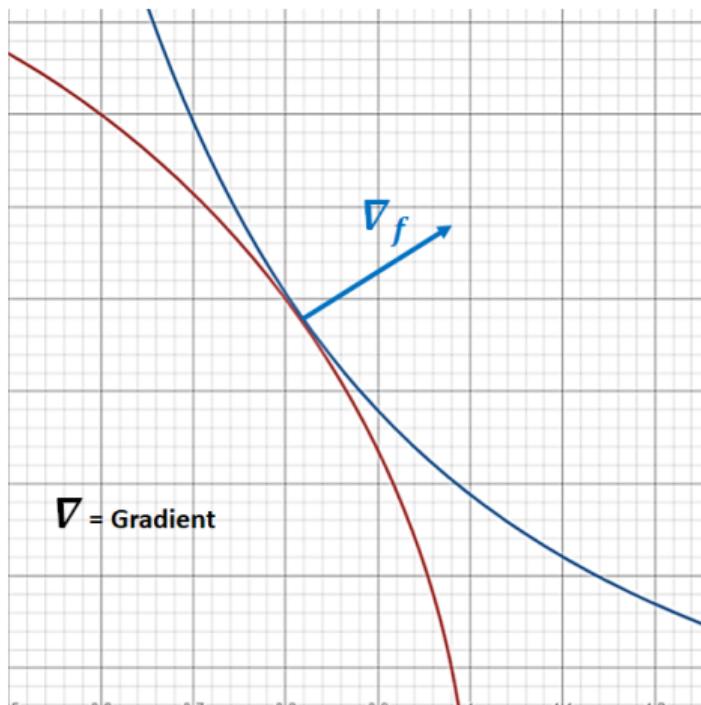
Optimizing function with constraint



The gradient of f or g evaluated at a point (x_0, y_0) always gives a vector perpendicular to the contour line passing through that point (as there is no change in value along contour line).

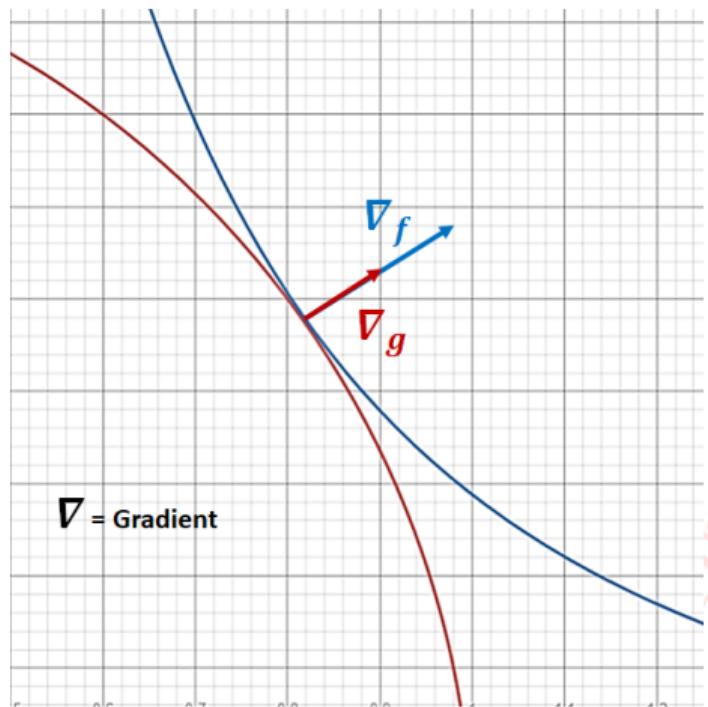
Function Optimization

Optimizing function with constraint



- The gradient of f evaluated at a point (x_0, y_0) always gives a vector perpendicular to the contour line passing through that point (as there is no change in value along contour line).
- When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.

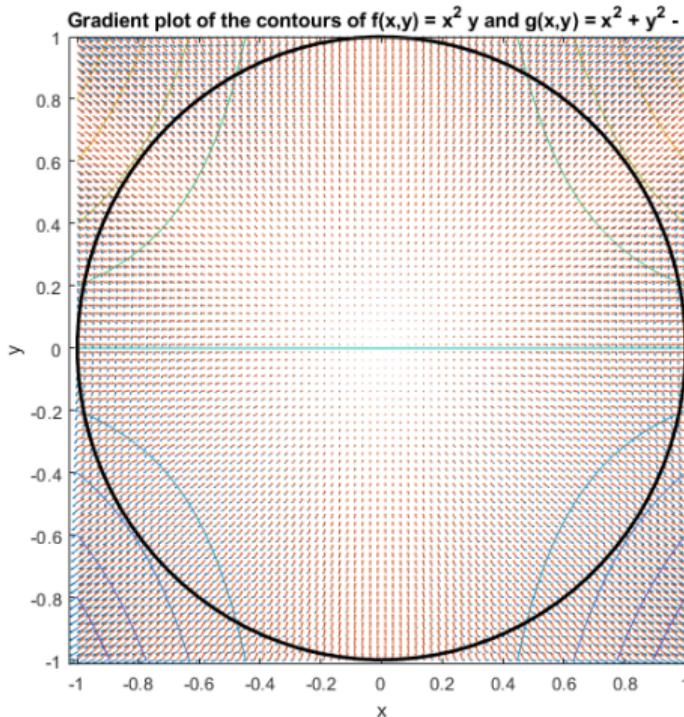
Optimizing function with constraint



- The fact that contour lines are tangent tells us nothing about the magnitude of each of these gradient vectors, but that's okay. When two vectors point in the same direction, it means we can multiply one by some constant to get the other.
- Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

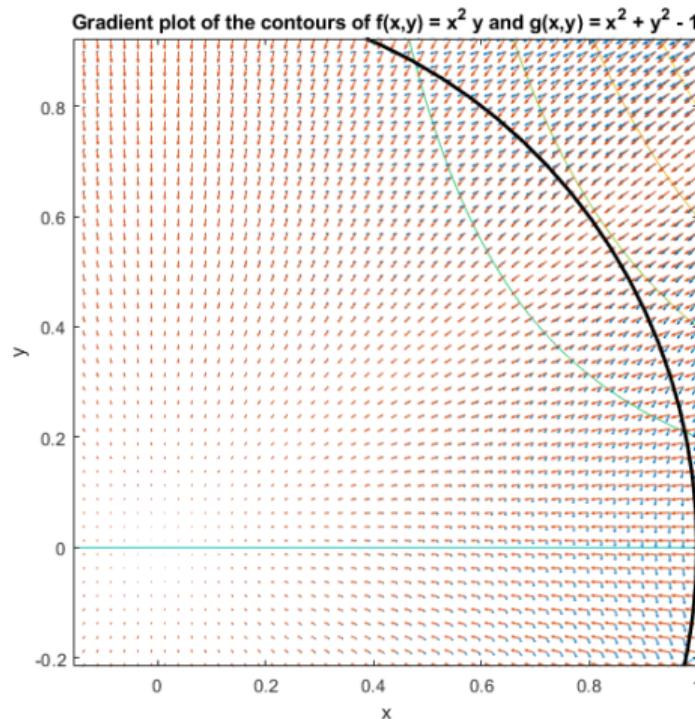
λ = Lagrange multiplier
 $f(x, y)$ = Function
 $g(x, y)$ = Constraint



- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
- Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

Ocular Proof

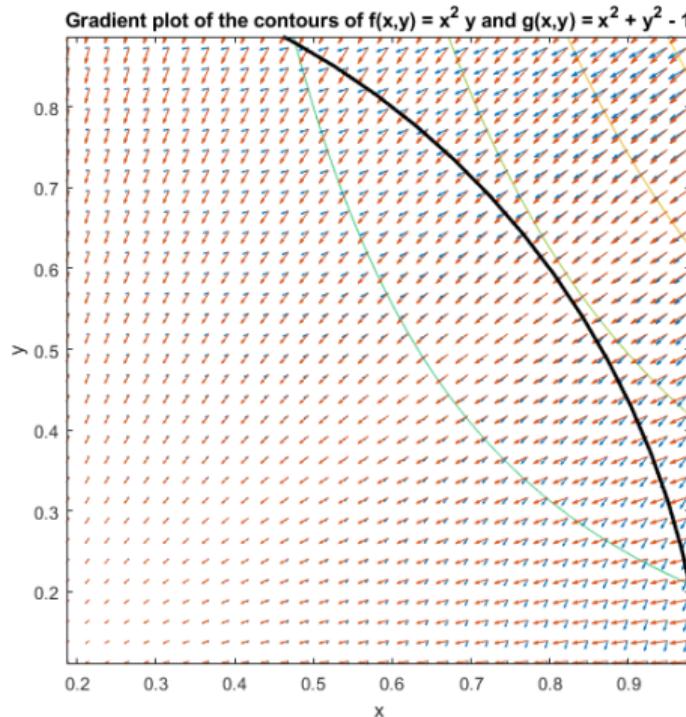


- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
- Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

Function Optimization

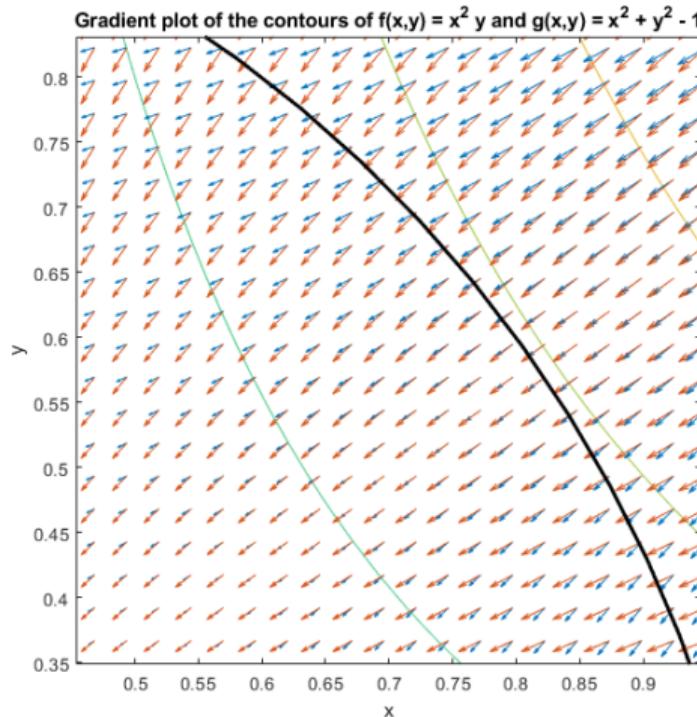
Ocular Proof



- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
- Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

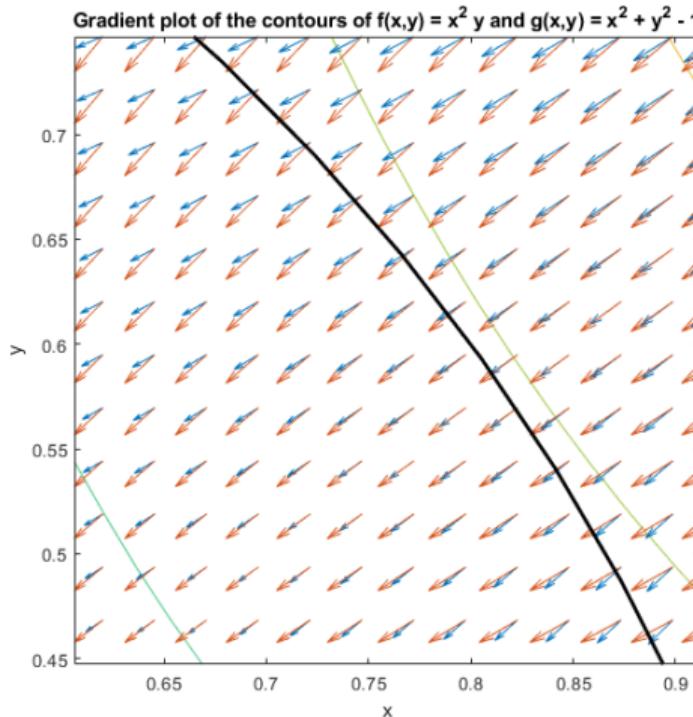
Ocular Proof



- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
 - Since this tangency means their gradient vectors align:

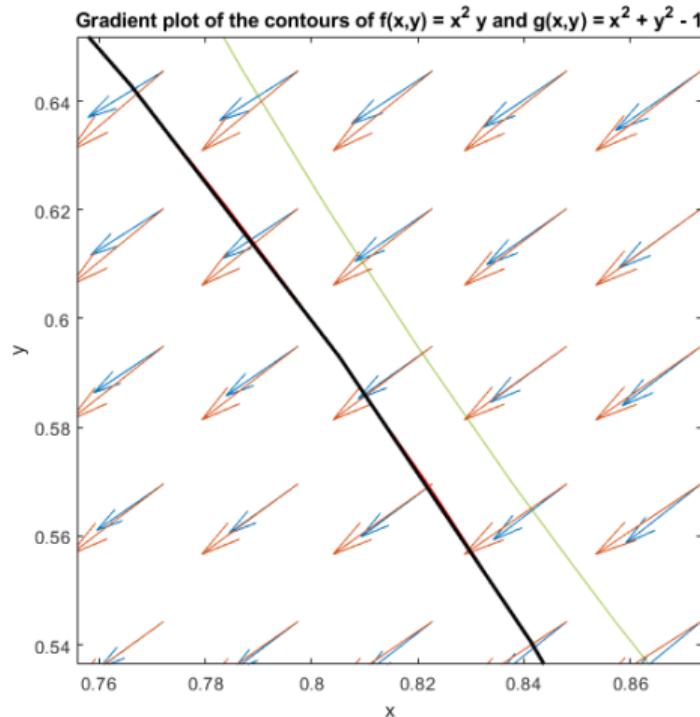
$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

Ocular Proof



- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
 - Since this tangency means their gradient vectors align:

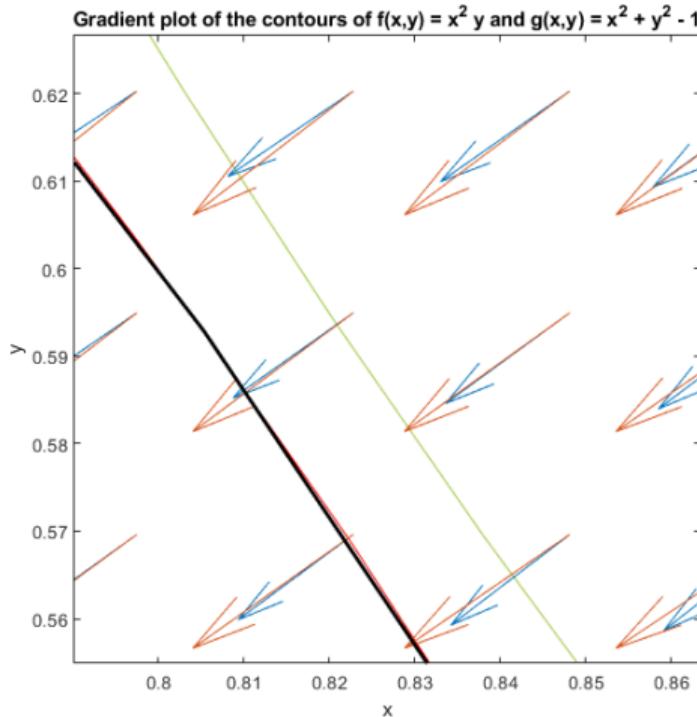
$$\nabla f(x, y) = \lambda \nabla g(x, y)$$



- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
- Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

Ocular Proof



A Khan

- **Ocular Proof:** When the contour lines of two functions f and g are tangent, their gradient vectors are parallel.
 - Since this tangency means their gradient vectors align:

$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

Optimizing function with constraint: Example Solution

$$\begin{aligned}L &= f(x, y) - \lambda g(x, y) \\&= x^2 y - \lambda [x^2 + y^2]\end{aligned}$$

To find max. take derivative. First partial derivative w.r.t "x":

$$\frac{\partial L}{\partial x} = 2xy - \lambda 2x \quad (31)$$
$$y = \lambda$$

Partial derivative w.r.t "y":

$$\frac{\partial L}{\partial y} = x^2 - \lambda 2y \quad (32)$$
$$x^2 = \lambda 2y$$
$$x^2 = 2\lambda^2 \text{ (as } y = \lambda\text{)}$$
$$x = \pm\sqrt{2}\lambda$$

Function Optimization

Optimizing function with constraint: Example Solution

Putting back values of “ x ” and “ y ” found from equations 31 and 32 in the constraint equation:

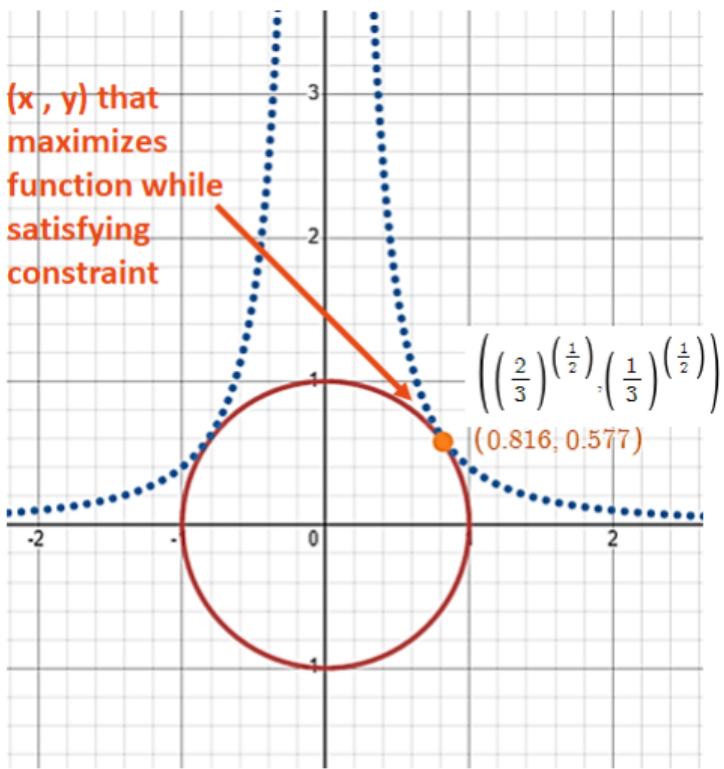
$$\begin{aligned}x^2 + y^2 &= 1 \\ [\sqrt{2}\lambda]^2 + \lambda^2 &= 1 \\ 3\lambda^2 &= 1 \\ \lambda &= \pm\sqrt{\frac{1}{3}}\end{aligned}\tag{33}$$

Put back value of λ in equations 31 and 32 to find value of x and y :

$$y = \pm\sqrt{\frac{1}{3}}\tag{34}$$

$$x = \pm\sqrt{2}\sqrt{\frac{1}{3}} = \pm\sqrt{\frac{2}{3}}\tag{35}$$

Optimizing function with constraint: Example Solution

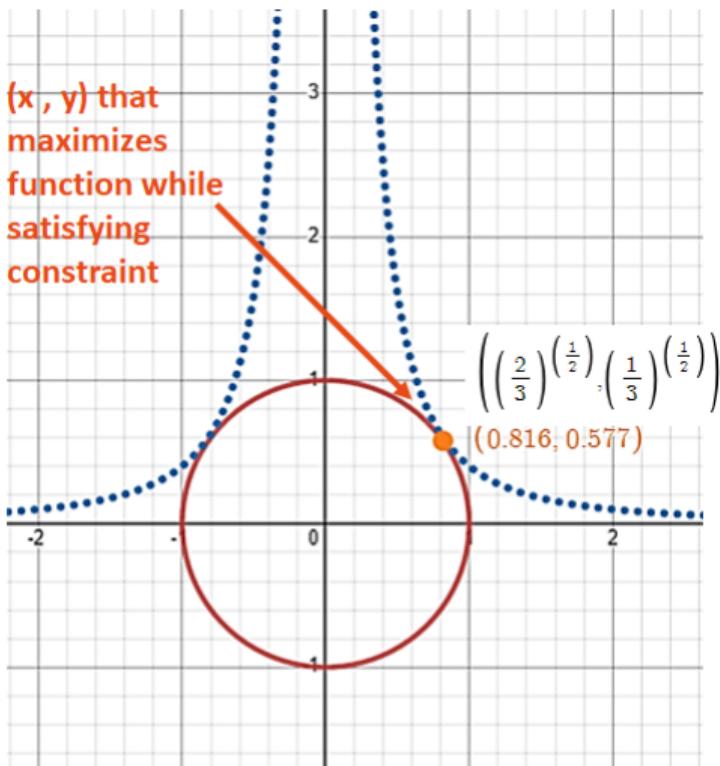


- From equations 34 and 35 , we know values of x and y . They make four possible pairs of (x, y) :

- 1** $(\sqrt{\frac{2}{3}}, \sqrt{\frac{1}{3}})$
 - 2** $(-\sqrt{\frac{2}{3}}, \sqrt{\frac{1}{3}})$
 - 3** $(\sqrt{\frac{2}{3}}, -\sqrt{\frac{1}{3}})$
 - 4** $(-\sqrt{\frac{2}{3}}, -\sqrt{\frac{1}{3}})$

Function Optimization

Optimizing function with constraint: Example Solution



- From equations 34 and 35 , we know values of x and y . They make four possible pairs of (x, y) :

- ① $(\sqrt{\frac{2}{3}}, \sqrt{\frac{1}{3}})$
- ② $(-\sqrt{\frac{2}{3}}, \sqrt{\frac{1}{3}})$
- ③ $(\sqrt{\frac{2}{3}}, -\sqrt{\frac{1}{3}})$
- ④ $(-\sqrt{\frac{2}{3}}, -\sqrt{\frac{1}{3}})$

- Last two point make function (x^2y) negative (will not max. func.) and first two points gives same output and that is the maximum value function can achieve while satisfying constraint (refer image on the left).