

- Midterm is Friday - Monday (10 days)
- Look at data, try to get some intuition, create some new and interesting features
- Singular value decomposition
- Will do well on the midterm if you find features that other people don't
- And the tuning parameters there of
- These are super powerful tools

- There is a way of generating multiple classifiers and then aggregating it into one mega classifier

- Midterm is same for both classes
- 80% of grade on report
- 20% proportional to the score you get
 - o Not necessarily your rank on the leader board
- Want to reward high level understanding of the tools and methods
- Have a good workflow and process
- Want every step to be well understood
- Each step should have some insight, and based on that you do the next step

- Is there every a reason not to use SVD?
 - o If you don't have a lot of features to begin with

- Deliverable 0 tonight
- Midterm Friday, start ASAP

Questions I have:

- Are midterm instructions posted?
- Will deliverable 1 be due during the midterm?
- Do we need to have a particular python version installed to use the provided start code?

```

File - C:\Users\Wolfs\PycharmProjects\cs506iec20211013\facial_recog.py
1 import numpy as np
2 from PIL import Image
3 import seaborn as sns
4 from sklearn.svm import SVC
5 import matplotlib.pyplot as plt
6 from sklearn.decomposition import PCA
7 from sklearn.pipeline import make_pipeline
8 from sklearn.metrics import confusion_matrix
9 from sklearn.datasets import fetch_lfw_people
10 from sklearn.ensemble import BaggingClassifier
11 from sklearn.model_selection import GridSearchCV,
12     train_test_split
13 # There are lots of datasets in sklearn that you would
14 # need to download
14 # Preprocess data to improve performance
15 # General data science flow for midterm
16
17 # This requires an older version of sklearn??
18 # It looks like he is using Python 3.7.4
19
20 sns.set()
21
22 # Get face data
23 faces = fetch_lfw_people(min_faces_per_person=60)
24 print(faces.target_names)
25
26 # plot face data
27 fig, ax = plt.subplots(3, 5)
28 for i, axi in enumerate(ax.flat):
29     axi.imshow(faces.images[i], cmap='bone')
30     axi.set(xticks=[], yticks=[], xlabel=faces.target_names[faces.target[i]])
31 plt.show()
32
33
34 # split train test set
35 Xtrain, Xtest, ytrain, ytest = train_test_split(faces.data, faces.target, random_state=42)
36
37 # blindly fit svm

```

Page 1 of 4

```

File - C:\Users\Wolfs\PycharmProjects\cs506iec20211013\facial_recog.py
38 svc = SVC(kernel='rbf', class_weight='balanced', C=10
39     , gamma=0.001)
40 """
41 #fit model
42 model = svc.fit(Xtrain, ytrain)
43 yfit = model.predict(Xtest)
44
45 fig, ax = plt.subplots(6, 6)
46 for i, axi in enumerate(ax.flat):
47     axi.imshow(Xtest[i].reshape(62, 47), cmap='bone')
48     axi.set(xticks=[], yticks=[], xlabel=faces.target_names[yfit[i]].split
49     ()[-1],
50             color='black' if yfit[i] == ytest[i]
51         ] else 'red')
51 fig.suptitle('Predicted Names; Incorrect Labels in Red',
52             size=14)
52 plt.show()
53
54 mat = confusion_matrix(ytest, yfit)
55 print(mat)
56 sns.heatmap(mat.T, square=True, annot=True, fmt='d',
57             cbar=False,
58             xticklabels=faces.target_names,
59             yticklabels=faces.target_names)
59 plt.xlabel('true label')
60 plt.ylabel('predicted label')
61 plt.show()
62 """
63
64 # Model just learned to predict bush for facial
64 # recognition
65 # look at confusion matrix to see what was predicted
65 # vs what it wanted
66
67 #####
68 #####
69 #####
70
71 # look at singular values

```

Page 2 of 4

```

File - C:\Users\Wolfs\PycharmProjects\cs509iec20211013\facial_recog.py
72 _, s, _ = np.linalg.svd(Xtrain, full_matrices=False)
73 plt.plot(range(1,len(s)+1),s)
74 plt.title("Singular Values")
75 plt.show()
76
77 # extract principal components
78 pca = PCA(n_components=150, whiten=True)
79 svc = SVC(kernel='rbf', class_weight='balanced')
80 # Make pipeline is worth looking up and see how it
ties toge
81 svcpca = make_pipeline(pca, svc)
82
83 param_grid = {'svc__C': [1, 5, 10, 50],
84                 'svc__gamma': [0.0001, 0.0005, 0.001, 0
     .005]}
85 grid = GridSearchCV(svcpca, param_grid)
86
87 # fit model
88 grid.fit(Xtrain, ytrain)
89 print(grid.best_params_)
90
91 model = grid.best_estimator_
92 yfit = model.predict(xtest)
93
94 fig, ax = plt.subplots(6, 6)
95 for i, axi in enumerate(ax.flat):
96     axi.imshow(xtest[i].reshape(62, 47), cmap='bone')
97     axi.set(xticks=[], yticks[])
98     axi.set_ylabel(faces.target_names[yfit[i]].split
()[-1],
99                   color='black' if yfit[i] == ytest[
i] else 'red')
100 fig.suptitle('Predicted Names; Incorrect Labels in
Red', size=14)
101 plt.show()
102
103 mat = confusion_matrix(ytest, yfit)
104 print(mat)
105 sns.heatmap(mat.T, square=True, annot=True, fmt='d',
cbar=False,
106                 xticklabels=faces.target_names,

```

Page 3 of 4

```

107                 yticklabels=faces.target_names)
108 plt.xlabel('true label')
109 plt.ylabel('predicted label')
110 plt.show()
111
112
113 # # optimal weights from before
114 # pca = PCA(n_components=20, whiten=True)
115 # svc = SVC(kernel='rbf', class_weight='balanced', C=
5, gamma=0.001)
116 # svcpca = make_pipeline(pca, svc)
117 # bagging = BaggingClassifier(svcpca, n_estimators=
100)
118 # model = bagging.fit(Xtrain, ytrain)
119 # yfit = model.predict(xtest)
120
121 # fig, ax = plt.subplots(6, 6)
122 # for i, axi in enumerate(ax.flat):
123 #     axi.imshow(xtest[i].reshape(62, 47), cmap='bone
')
124 #     axi.set(xticks=[], yticks[])
125 #     axi.set_ylabel(faces.target_names[yfit[i]].split
()[-1],
126 #                   color='black' if yfit[i] ==
ytest[i] else 'red')
127 # fig.suptitle('Predicted Names; Incorrect Labels in
Red', size=10)
128 # plt.show()
129
130 # mat = confusion_matrix(ytest, yfit)
131 # print(mat)
132 # sns.heatmap(mat.T, square=True, annot=True, fmt='d
',
cbar=False,
133 #                 xticklabels=faces.target_names,
134 #                 yticklabels=faces.target_names)
135 # plt.xlabel('true label')
136 # plt.ylabel('predicted label')
137 # plt.show()
138

```

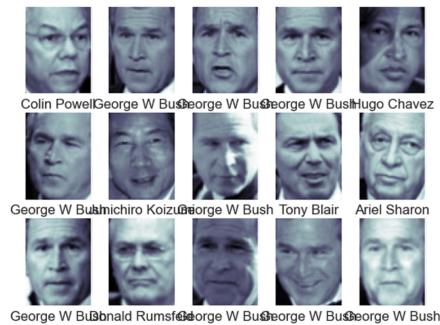
Page 4 of 4

```
1 cycler==0.10.0
2 joblib==0.14.0
3 kiwisolver==1.1.0
4 matplotlib==3.1.1
5 numpy==1.17.3
6 pandas==0.25.3
7 Pillow==6.2.1
8 pyParsing==2.4.4
9 python-dateutil==2.8.1
10 pytz==2019.3
11 scikit-learn==0.21.3
12 scipy==1.3.1
13 seaborn==0.9.0
14 six==1.13.0
15 sklearn==0.0
16
```

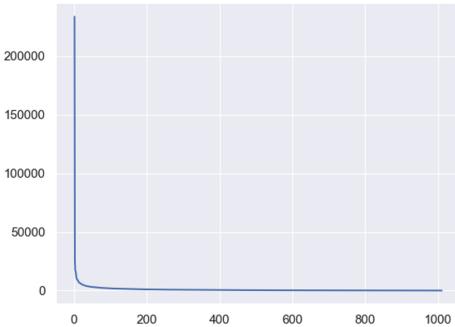
C:\Users\Wolfs\anaconda3\envs\cs506lec20211013\python.exe C:/Users/Wolfs/PycharmProjects/cs506lec20211013/facial_recog.py

```
['Ariel Sharon' 'Colin Powell' 'Donald Rumsfeld' 'George W Bush'
 'Gerhard Schroeder' 'Hugo Chavez' 'Junichiro Koizumi' 'Tony Blair']
{'svc_C': 5, 'svc_gamma': 0.001}
[[ 15  0  2  0  0  0  0  0]
 [ 2  61  2  3  0  0  0  0]
 [ 1  2  26  0  1  0  0  1]
 [ 2  9  4 104  1  1  2  3]
 [ 0  0  1  1 18  1  0  2]
 [ 1  1  0  0  0 16  0  2]
 [ 0  0  0  0  0  0 12  0]
 [ 1  0  0  0  1  0  0 40]]
```

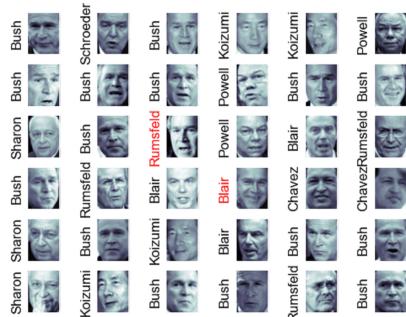
Process finished with exit code 0



Singular Values



Predicted Names; Incorrect Labels in Red



Ariel Sharon	13	2	1	2	0	1	0	1
Colin Powell	0	61	2	9	0	1	0	0
Donald Rumsfeld	2	2	26	4	1	0	0	0
George W Bush	0	3	0	104	1	0	0	0
Gerhard Schroeder	0	0	1	1	18	0	0	1
Hugo Chavez	0	0	0	1	1	16	0	0
Junichiro Koizumi	0	0	0	2	0	0	12	0
Tony Blair	0	0	1	3	2	2	0	40

sharon
owell
nsfeld
nsfield
Bush
soder
avez
izumi
r Blair

20211013 Naïve Bayes and SVM

Wednesday, October 13, 2021 4:57 PM

Naïve Bayes and SVM

Bayes Classifier

- A probabilistic framework for solving classification problems

- Conditional Probability: $P(C | A) = \frac{P(A, C)}{P(A)}$

$$P(A | C) = \frac{P(A, C)}{P(C)}$$

- Bayes theorem: $P(C | A) = \frac{P(A | C)P(C)}{P(A)}$

Example of Bayes Theorem

- Given:
 - A doctor knows that meningitis causes stiff neck 50% of the time
 - Prior probability of any patient having meningitis is 1/50,000
 - Prior probability of any patient having stiff neck is 1/20
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

- Probability of stiff neck with meningitis is high, but if you have a stiff neck that doesn't really mean high probability of having meningitis

Bayesian Classifiers

- Consider each attribute and class label as random variables
- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate $P(C | A_1, A_2, \dots, A_n)$ directly from data?
 - Have to know the probability of A1, A2,... and so on
 - If they are independent that is a little easier
 - Need a prior probability of being in C as well

Bayesian Classifiers

- Approach:
 - compute the posterior probability $P(C | A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem
 - Choose value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
 - Equivalent to choosing value of C that maximizes $P(A_1, A_2, \dots, A_n | C) P(C)$
- How to estimate $P(A_1, A_2, \dots, A_n | C)$?

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

Naïve Bayes Classifier

- Assume independence among attributes A_i when class is given:
 - $P(A_1, A_2, \dots, A_n | C) = P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$
 - Can estimate $P(A_i | C_j)$ for all A_i and C_j .
 - New point is classified to C_j if $P(C_j) \prod P(A_i | C_j)$ is maximal.
- e.g. height, weight, bmi
 - BMI is a function of height and weight
 - These are NOT independent
- Assuming independence is a BIG assumption
- If a single one of these is 0 then you loose a lot of information (multiplying all the probabilities together)

How to Estimate Probabilities from Data?

- Class: $P(C) = N_c/N$

– e.g., $P(\text{No}) = 7/10$,
 $P(\text{Yes}) = 3/10$

- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_c$$

– where $|A_{ik}|$ is number of instances having attribute A_i and belongs to class C_k

- Examples:

$$\begin{aligned} P(\text{Status=Married} | \text{No}) &= 4/7 \\ P(\text{Refund=Yes} | \text{Yes}) &= 0 \end{aligned}$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

How to Estimate Probabilities from Data?

- For continuous attributes:
 - Discretize the range into bins
 - one ordinal attribute per bin
 - violates independence assumption
 - Two-way split: $(A < v)$ or $(A > v)$
 - choose only one of the two splits as new attribute
 - Probability density estimation:
 - Assume attribute follows a normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
 - Once probability distribution is known, can use it to estimate the conditional probability $P(A_i | c)$

-

How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

– One for each (A_i, c_j) pair

- For (Income, Class=No):

– If Class=No

- sample mean = 110
- sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Example of Naïve Bayes Classifier

Given a Test Record:

$$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$$

naive Bayes Classifier:

$P(\text{Refund}=\text{Yes} \text{No}) = 3/7$ $P(\text{Refund}=\text{No} \text{No}) = 4/7$ $P(\text{Refund}=\text{Yes} \text{Yes}) = 0$ $P(\text{Refund}=\text{No} \text{Yes}) = 1$ $P(\text{Marital Status}=\text{Single} \text{No}) = 2/7$ $P(\text{Marital Status}=\text{Divorced} \text{No}) = 1/7$ $P(\text{Marital Status}=\text{Married} \text{No}) = 4/7$ $P(\text{Marital Status}=\text{Single} \text{Yes}) = 2/7$ $P(\text{Marital Status}=\text{Divorced} \text{Yes}) = 1/7$

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \times P(\text{Married}|\text{Class}=\text{No}) \times P(\text{Income}=120\text{K}|\text{Class}=\text{No})$
 $= 4/7 \times 4/7 \times 0.0072 = 0.0024$

- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes}) \times P(\text{Married}|\text{Class}=\text{Yes}) \times P(\text{Income}=120\text{K}|\text{Class}=\text{Yes})$

$$\begin{aligned} P(\text{Marital Status}=\text{Married}|\text{No}) &= 4/7 \\ P(\text{Marital Status}=\text{Single}|\text{Yes}) &= 2/7 \\ P(\text{Marital Status}=\text{Divorced}|\text{Yes}) &= 1/7 \\ P(\text{Marital Status}=\text{Married}|\text{Yes}) &= 0 \end{aligned}$$

For taxable income:

If class=No:	sample mean=110
	sample variance=2975
If class=Yes:	sample mean=90
	sample variance=25

- $P(X|\text{Class}=\text{Yes}) = P(\text{Retuna}=\text{No} | \text{Class}=\text{Yes})$
 $\times P(\text{Married} | \text{Class}=\text{Yes})$
 $\times P(\text{Income}=120K | \text{Class}=\text{Yes})$
 $= 1 \times 0 \times 1.2 \times 10^{-9} = 0$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
 $\Rightarrow \text{Class} = \text{No}$

- =0 is an issue

Naïve Bayes Classifier

- If one of the conditional probability is zero, then the entire expression becomes zero
- Probability estimation:

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c} \quad c: \text{number of classes}$$

$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c} \quad p: \text{prior probability}$$

$$\text{m - estimate : } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m} \quad m: \text{parameter}$$

- Look at 1+ that so that it is never 0
- add c based on how many features you have so that you have a proper probability
- A probability of 0 is you don't have a particular

example in your dataset

- But this could be because you don't have enough data
- If you have A LOT of data and you still have a 0, then the 0 might be real valid estimate
- Can use these different transforms in this case, normalizing so they add up to proper probabilities

Example of Naïve Bayes Classifier

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

A: attributes

M: mammals

N: non-mammals

$$P(A|M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A|N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A|M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A|N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

$$P(A|M)P(M) > P(A|N)P(N)$$

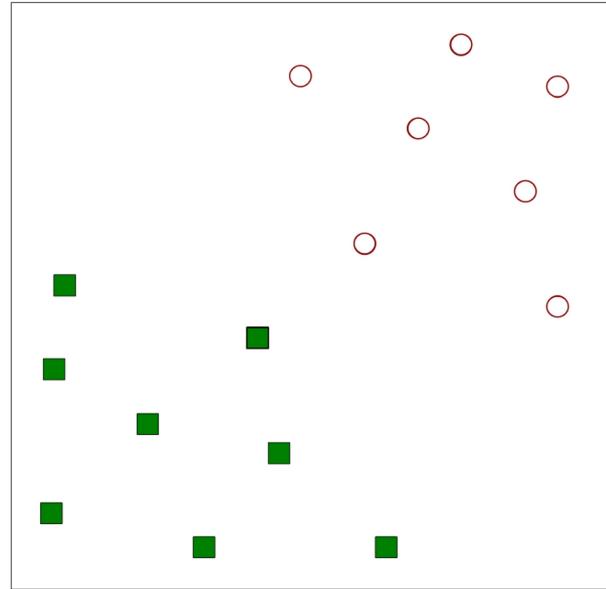
=> Mammals

- What is the probability it's mammal given these attributes
- What is the probability it's NOT mammal given these attributes
- Assign to whatever has the higher probability
- Assume independence to make the problem easier

Naïve Bayes (Summary)

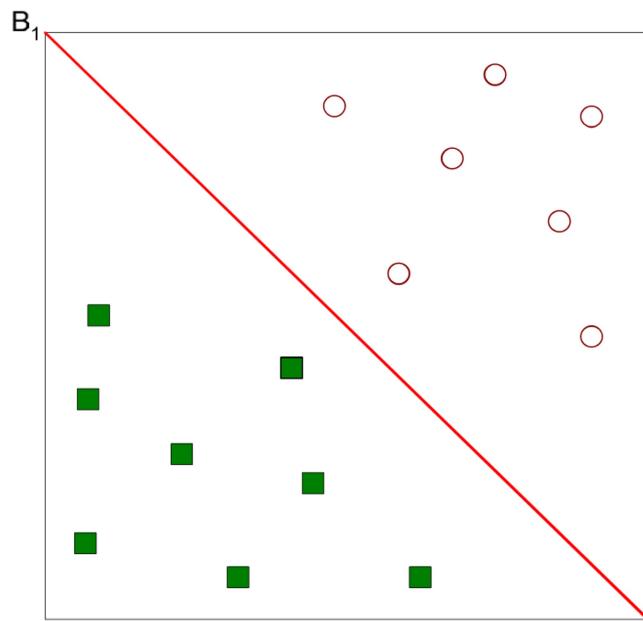
- Robust to isolated noise points
 - Handle missing values by ignoring the instance during probability estimate calculations
 - Robust to irrelevant attributes
 - Independence assumption may not hold for some attributes
 - Use other techniques such as Bayesian Belief Networks (BBN)
- Called Naïve because it's assuming independence

Support Vector Machines



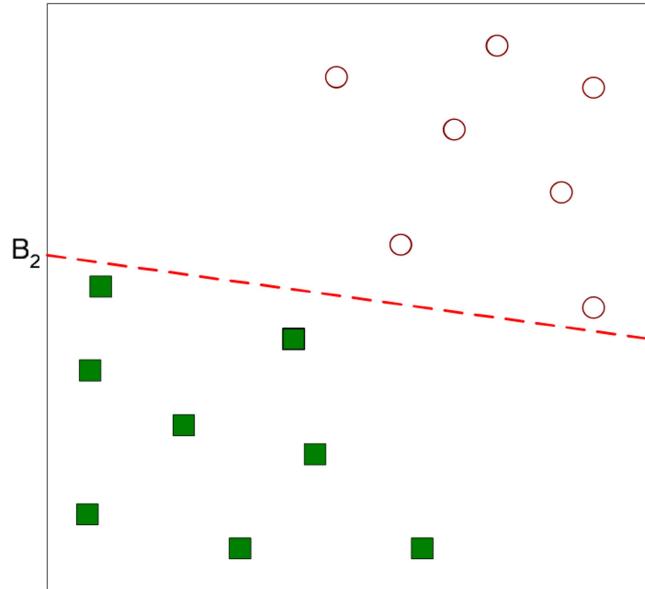
- Find a linear hyperplane (decision boundary) that will separate the data

Support Vector Machines



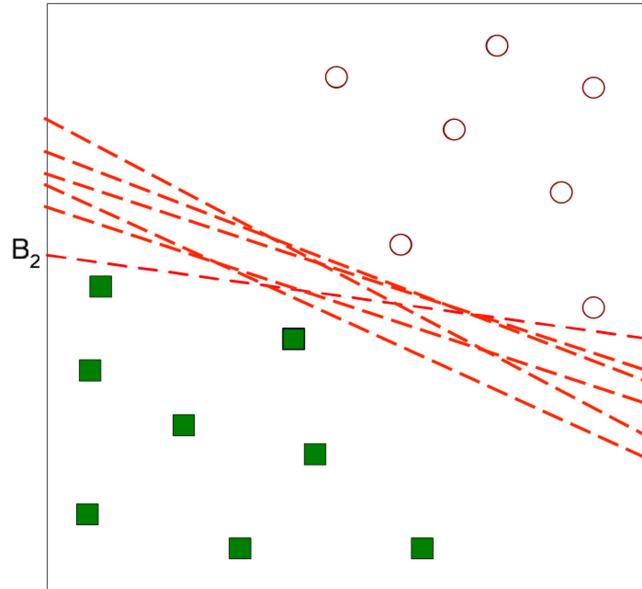
- One Possible Solution

Support Vector Machines



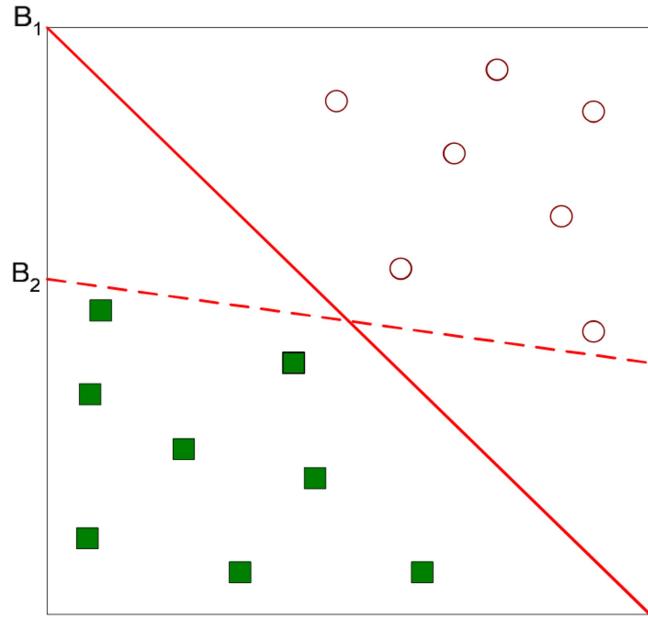
- Another possible solution

Support Vector Machines



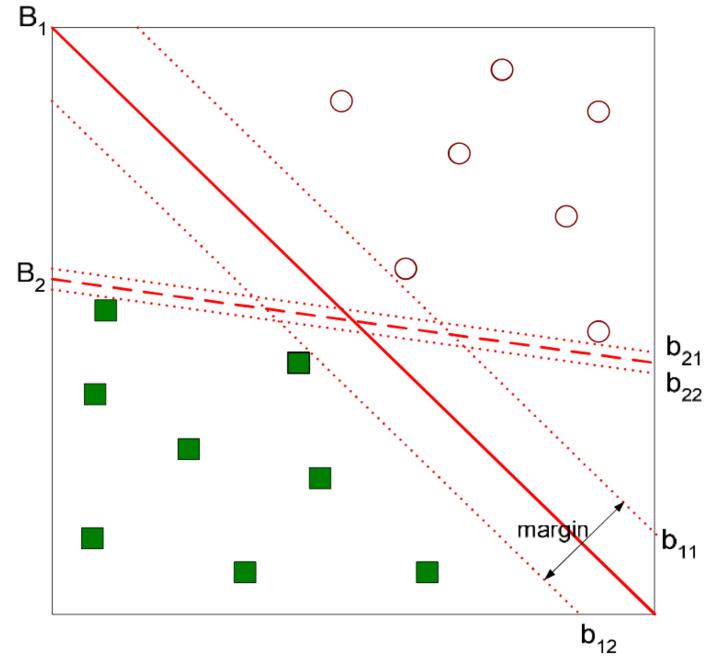
- Other possible solutions

Support Vector Machines



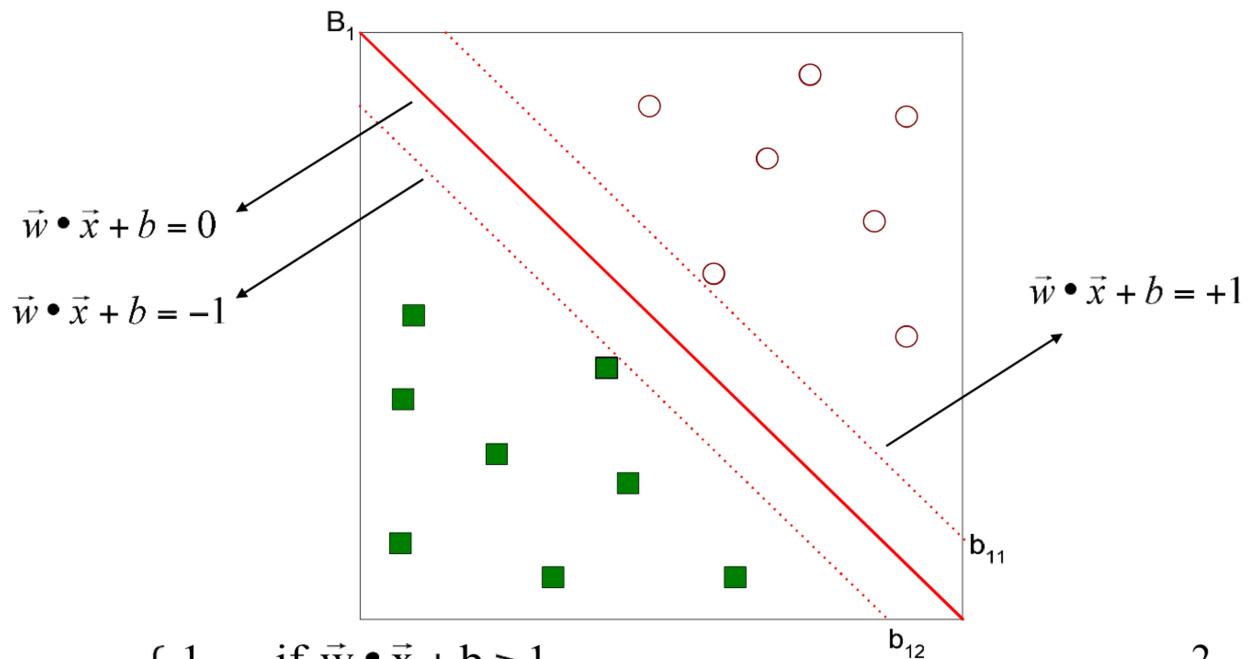
- Which one is better? B_1 or B_2 ?
- How do you define better?

Support Vector Machines



- Find hyperplane **maximizes** the margin => B_1 is better than B_2

Support Vector Machines



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

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

Support Vector Machines

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

- We want to maximize:

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

- Which is equivalent to minimizing:

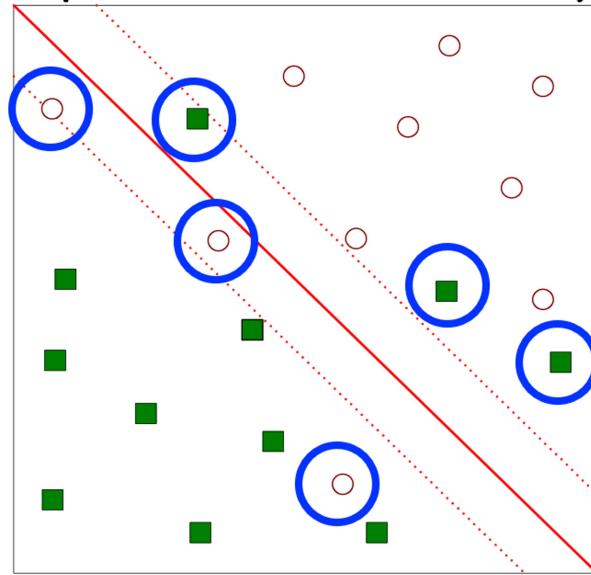
- But subjected to the following constraints:

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

- This is a constrained optimization problem
 - Numerical approaches to solve it (e.g., quadratic programming)

Support Vector Machines

- What if the problem is not linearly separable?



Support Vector Machines

- What if the problem is not linearly separable?

– Introduce slack variables

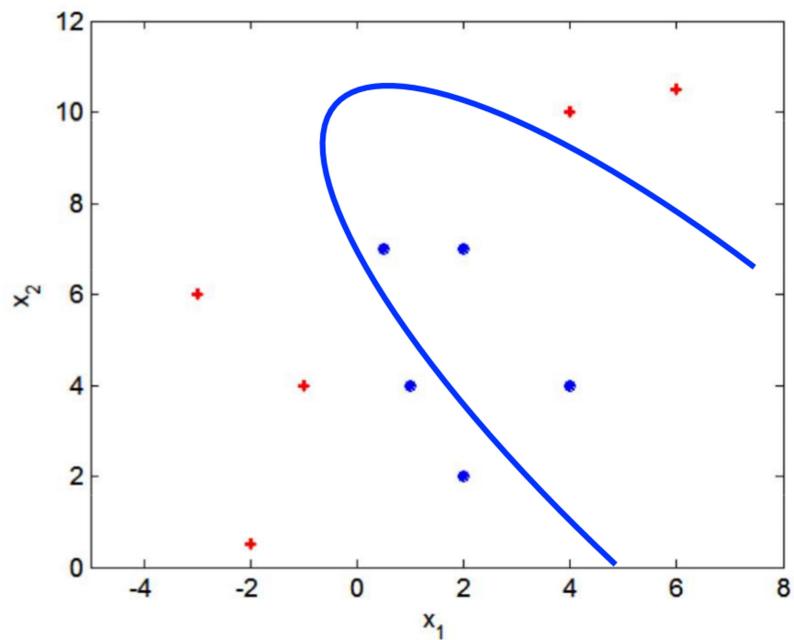
- Need to minimize: $L(w) = \frac{\|\vec{w}\|^2}{2} + C \left(\sum_{i=1}^N \xi_i \right)$

- Subject to:

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

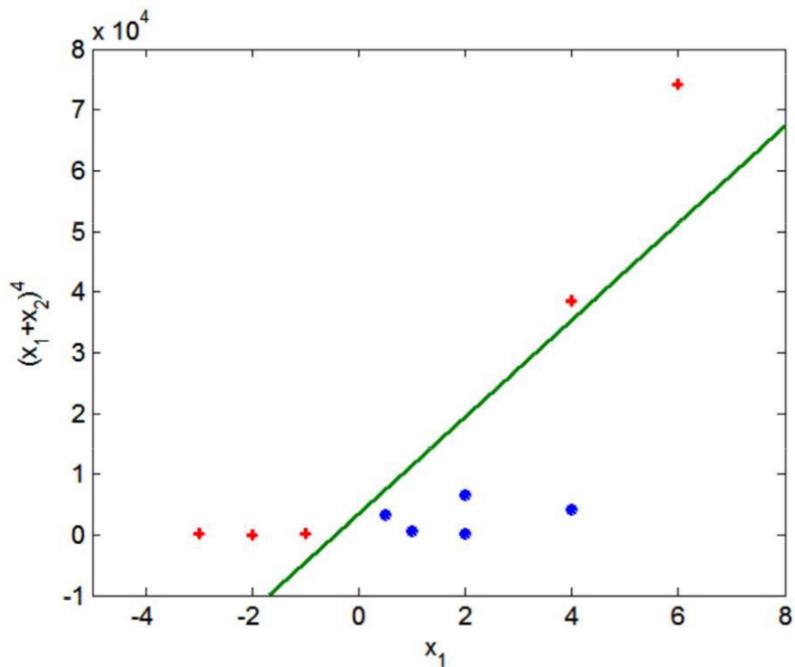
Nonlinear Support Vector Machines

- What if decision boundary is not linear?



Nonlinear Support Vector Machines

- Transform data into higher dimensional space



Kernel Trick!

- We do not need to actually map explicitly each point to a high dimensional space!!
- We just need to have a function that computes the similarity (distance) in the mapped space given the points in the input space (without the need to do the mapping!!!)
- Kernel function!!! (unfortunate term 😞)

Kernel function

- Given x and y two vectors in \mathbb{R}^d assume that we want to map them in a higher dimensional space \mathbb{R}^n using a function ϕ
- We can define a function $K(x,y): \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ that computes directly the distance of $\phi(x)$ and $\phi(y)$!
No need to compute the actual ϕ !!
- Actually, the Kernel function K should give the same results as the inner (dot) product between of $\phi(x)$ and $\phi(y)$ (simulates a dot product in the “mapped” space)

$$K(x,y) = \langle \phi(x), \phi(y) \rangle_n$$

SVM with Kernels

- Using the kernel trick, we can compute the linear separator in the mapped space without the need to do the actual mapping!! just use the Kernel function when you need to compute the inner product between two “mapped” vectors.

Kernel options

- Many different ways to define kernels
- Most popular kernels:
 - Linear: $K(x,y) = (x \bullet y)$
 - Polynomial: $K(x,y) = (x \bullet y)^d$ or $K(x,y) = (x \bullet y+1)^d$
 - RBF: $K(x,y) = \exp(-\gamma ||x-y||^2)$

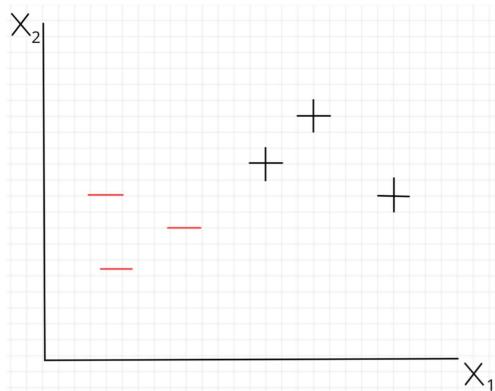
In practice, RBF works well in most cases (adds small bumps around the “mapped” points)

- If you don't have access to GIT
 - o Keep harassing everyone with cc'ing to address the issue

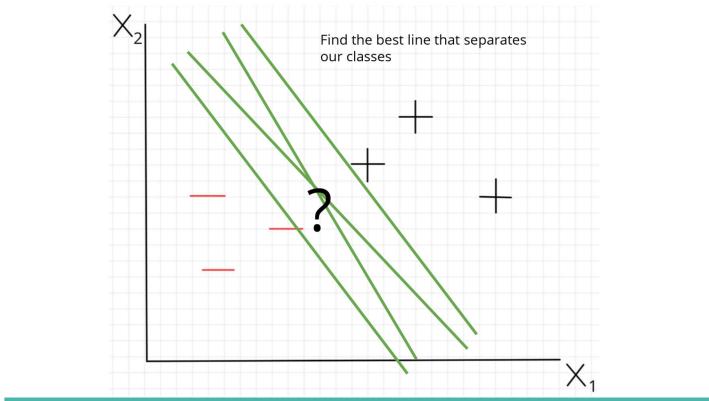
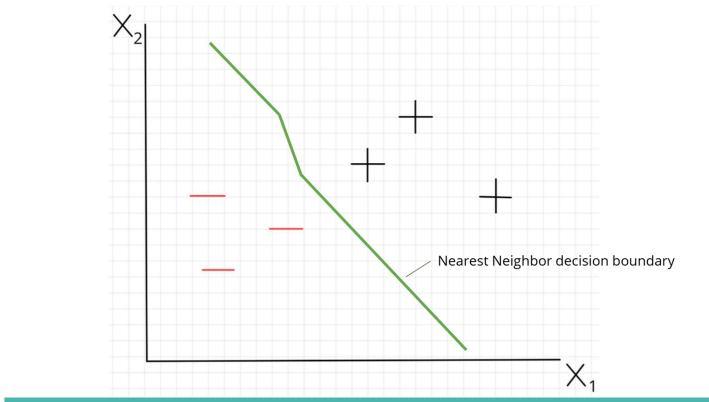
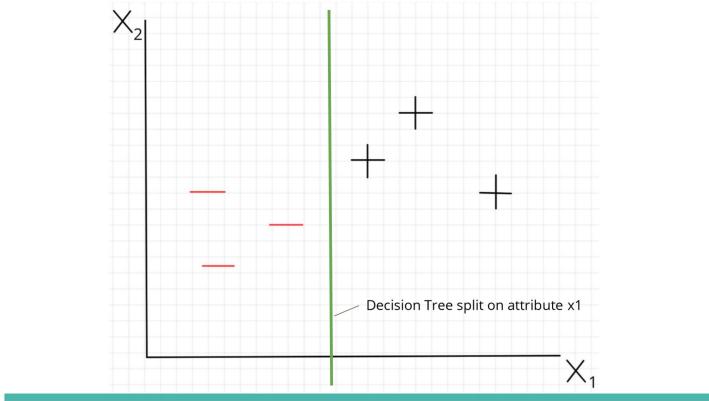
Support Vector Machines

Boston University CS 506 - Lance Galletti

- Coming back to the example we had at the beginning of class

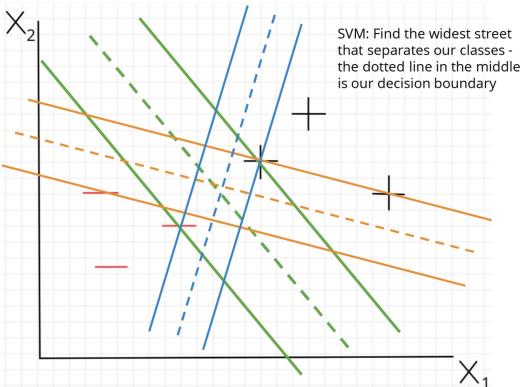


- We've seen a lot of classification techniques in the last day

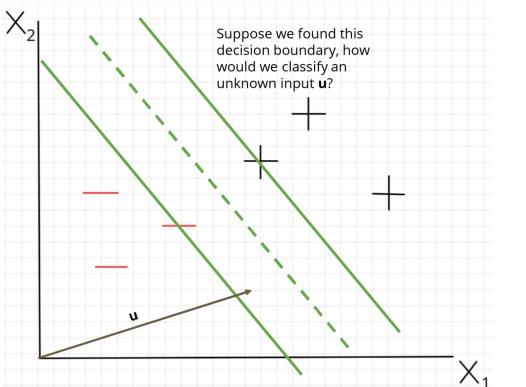


- There are a lot of possible lines, so what do we

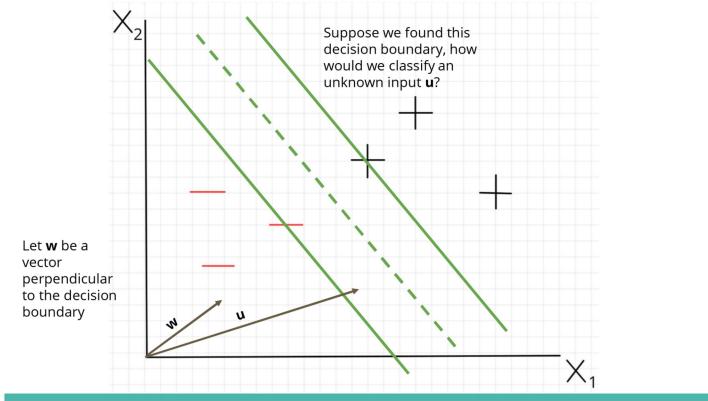
mean by the best line that separates our data?



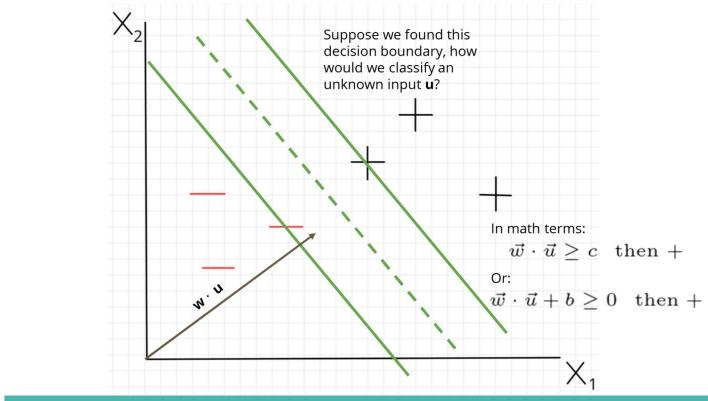
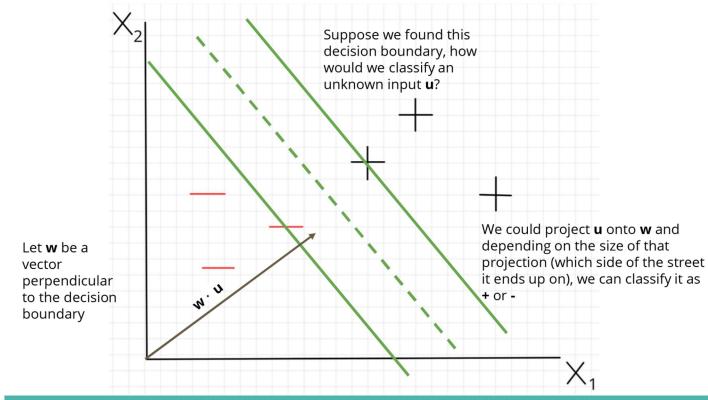
- SVM: tries to maximize the separation between the classes
- In this case, it will be the green one
- How do we find this street?



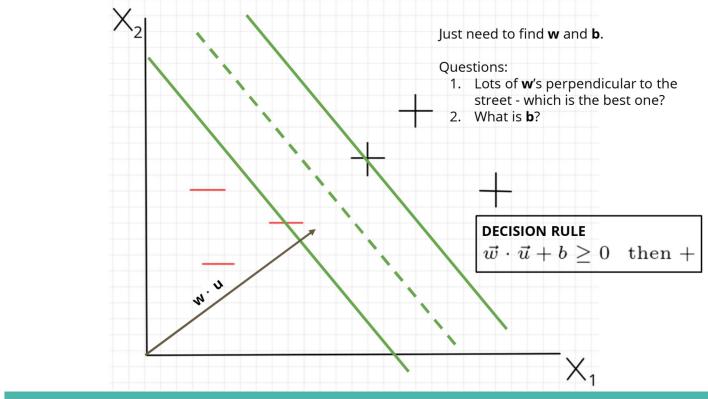
- Once we have the widest street, how do you classify unseen record
- What side of the street is u on?
 - o Find the vector that is perpendicular of decision boundary
 - o



- $w \cdot u$ (dot)



- bias term "b"



- We can fully characterize this street if we have the right w and b

How to find the widest street

We want our samples to lie beyond the street. That is:

$$\vec{w} \cdot \vec{x}_+ + b \geq 1$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1$$

Note: for an unknown u , we can have

$$-1 < \vec{w} \cdot \vec{u} + b < 1$$

-
- Whatever w we find must satisfy this condition
 - We don't WANT any examples to be inside the street
 - Want some space between positive and negative classes
 - -1 & 1 is just arbitrary
 - Want to find w and b that maximizes the width of that street
 - Unknowns can fall inside the street later
 - This is pretty complicated

How to find the widest street

Let's introduce a variable

$$y_i = \begin{cases} +1 & \text{if } x_i \text{ is a + sample} \\ -1 & \text{if } x_i \text{ is a - sample} \end{cases}$$

How to find the widest street

Let's introduce a variable

$$y_i = \begin{cases} +1 & \text{if } x_i \text{ is a + sample} \\ -1 & \text{if } x_i \text{ is a - sample} \end{cases}$$

Note: this is effectively the class label of x_i

How to find the widest street

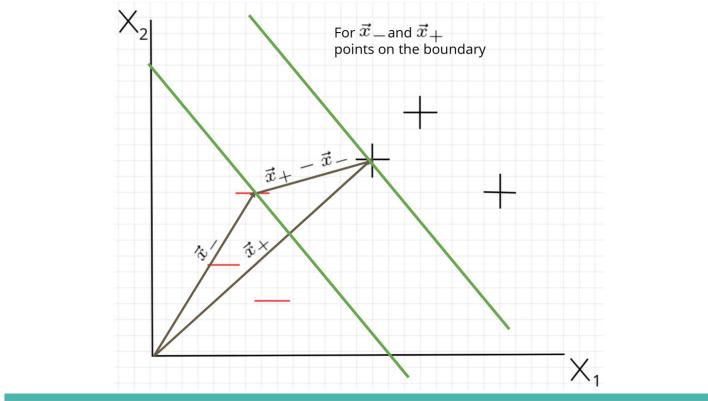
If we multiply our sample decision rules by this new variable:

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$

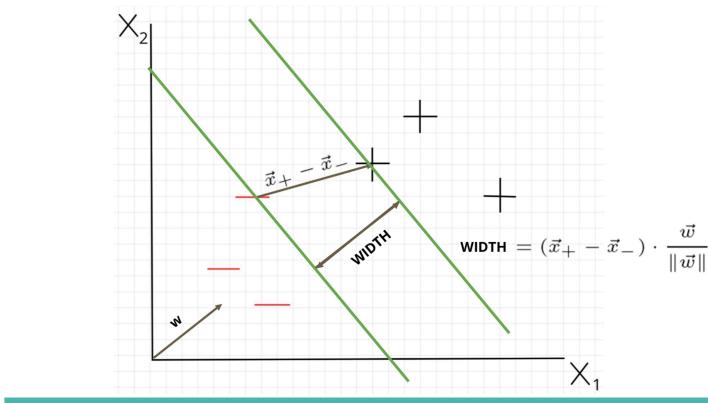
Meaning, for x_i on the decision boundary, we want:

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$

- If we multiply both by y_i , then we get these equations
- Now we have both constraints inside a single constraint
- Now we have equation of what to do with points EXACTLY on the side of the street



- Now we need to find the width of the street
- x_+ minus x_-



- Now have the width of the street
- Two points can fully characterize this?

Are there multiple possible solutions? (especially with a small number of points)

How to find the widest street

We know that $\text{WIDTH} = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|}$ for \vec{x}_- and \vec{x}_+ points on the boundary

And, since they are on the boundary, we know that

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$

$$\text{Hence, WIDTH} = \frac{2}{\|\vec{w}\|}$$

(as an exercise, try to show this)

How to find the widest street

Goal is to maximize the width

$$\begin{aligned} \max\left(\frac{2}{\|\vec{w}\|}\right) &= \min(\|\vec{w}\|) \\ &= \min\left(\frac{1}{2} \|\vec{w}\|^2\right) \end{aligned}$$

Subject to:

$$y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$

How to find the widest street

Can use Lagrange multipliers to form a single expression to find the extremum of

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum_i \alpha_i [y_i(\vec{x}_i \cdot \vec{w} + b) - 1]$$

where α_i is 0 for \vec{x}_i not on the boundary.

Now we can take derivatives to find the extremum of L .

- Take derivates of L and partial derivative with respect to w and b

How to find the widest street

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_i \alpha_i y_i \vec{x}_i = 0$$

$$\Rightarrow \vec{w} = \sum_i \alpha_i y_i \vec{x}_i$$

Means \mathbf{w} is a linear sum of vectors in our sample/training set!

$$\frac{\partial L}{\partial b} = - \sum_i \alpha_i y_i = 0$$

$$\Rightarrow \sum_i \alpha_i y_i = 0$$

- \mathbf{w} just a linear combination of things in the dataset, points on the boundary between the two classes

How to find the widest street

Let's plug these values back into L to see what happens to L at its extremum

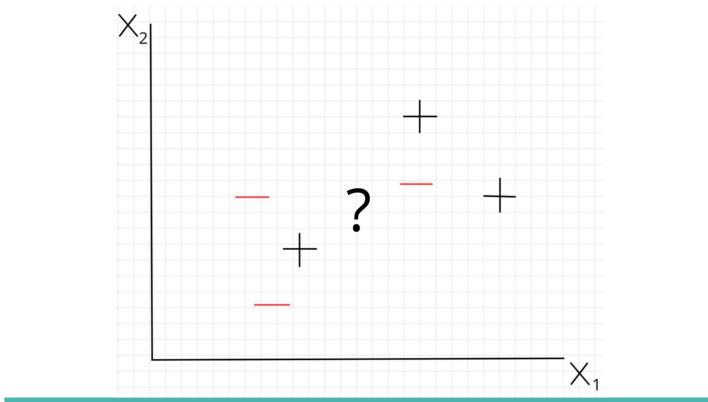
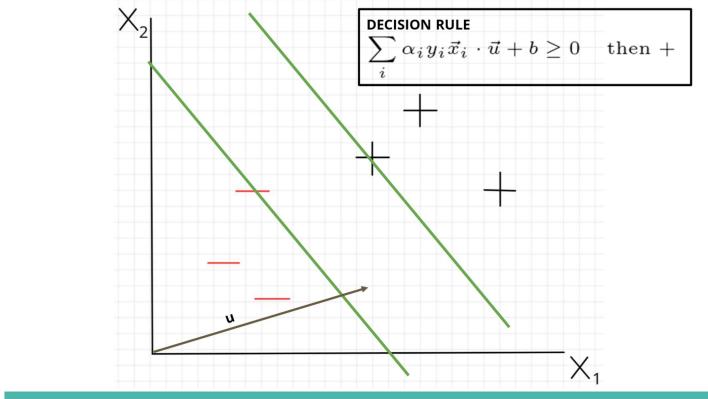
$$L = \frac{1}{2} \left(\sum_i \alpha_i y_i \vec{x}_i \right) \cdot \left(\sum_i \alpha_i y_i \vec{x}_i \right) - \left(\sum_i \alpha_i y_i \vec{x}_i \right) \cdot \left(\sum_i \alpha_i y_i \vec{x}_i \right) - \sum_j \cancel{\alpha_i y_i b} + \sum_i \alpha_i$$

Simplifying, we get:

$$L = \sum_i \alpha_i - \frac{1}{2} \left(\sum_i \alpha_i y_i \vec{x}_i \right) \cdot \left(\sum_i \alpha_i y_i \vec{x}_i \right)$$

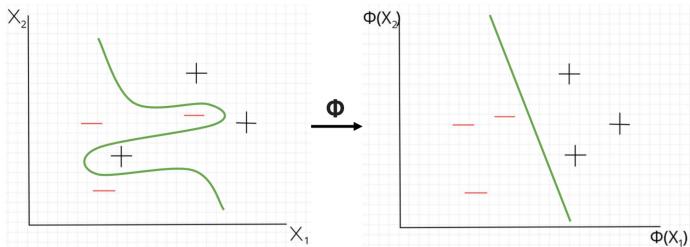
$$= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \cancel{\vec{x}_i \cdot \vec{x}_j}$$

- What happens to this optimization function as we get to the extreme?
- Linear combination of points in this dataset
- Fully characterized by inner products of points in the dataset



- What if our data is not linearly separable?
- DO TRANSFORMATION PHI and then do this
- There are a bunch of transformations that could achieve linearly separable
 - o infinite many, so lets define a function

When stuck - change perspective



But how to find Φ ?

Turns out we don't need to find or define a transformation Φ !

Looking back at L, since it depends only on the dot product of our input, we only need to define the dot product in our transformed space.

i.e. we only need to define

$$K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

Called a Kernel function. This is often referred to as the "kernel trick".

-
- You don't necessarily need to find the transform, you just need to know what the inner product looks like in that space
 - Just need to define the inner product for a potential candidate for the space

Example Kernel Functions

$$K(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j + 1)^n$$

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

DEMO
