

# Where Are We?

[Remove Watermark Now](#)

 pdfelement

# Where Are We?

Remove Watermark Now



# Where Are We?

Remove Watermark Now

We've looked some simple (but high performing) supervised learning ideas for estimating the class of individual documents (Naive Bayes), and for estimating proportions.



# Where Are We?

Remove Watermark Now

We've looked some simple (but high performing) **supervised learning** ideas for estimating the class of individual documents (Naive Bayes), and for estimating proportions.

Now want to cover powerful, commonly used techniques from **machine learning** that appear in social science research:



# Where Are We?

Remove Watermark Now

We've looked some simple (but high performing) **supervised learning** ideas for estimating the class of individual documents (Naive Bayes), and for estimating proportions.

Now want to cover powerful, commonly used techniques from **machine learning** that appear in social science research: **SVM**, KNN, CART etc.



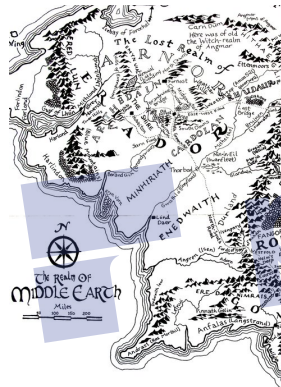
# Where Are We?

Remove Watermark Now

We've looked some simple (but high performing) **supervised learning** ideas for estimating the class of individual documents (Naive Bayes), and for estimating proportions.

Now want to cover powerful, commonly used techniques from **machine learning** that appear in social science research: **SVM**, KNN, CART etc.

These techniques involve some important decisions about the **bias-variance** tradeoff, and the use of **(cross) validation** in checking model performance and selecting the best model.



Remember...

[Remove Watermark Now](#)

 pdfelement

# Remember...

[Remove Watermark Now](#)

Unsupervised techniques:

pdfelement



# Remember...

[Remove Watermark Now](#)

Unsupervised techniques: learning  
(hidden or latent) structure in  
unlabeled data.

e.g. PCA of legislators's votes:

pdfelement

# Remember. . .

[Remove Watermark Now](#)

Unsupervised techniques: learning  
(hidden or latent) structure in  
unlabeled data.

e.g. PCA of legislators's votes: want to see  
how they are organized—



pdfelement

# Remember...

[Remove Watermark Now](#)

Unsupervised techniques: learning  
(hidden or latent) structure in  
unlabeled data.

e.g. PCA of legislators's votes: want to see  
how they are organized—by party? by  
ideology? by race?

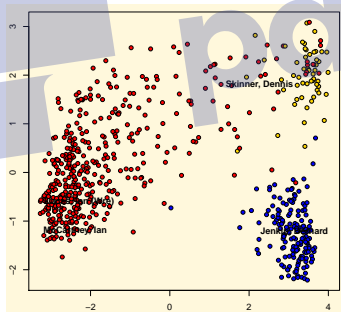
pdfelement

# Remember. . .

Remove Watermark Now

Unsupervised techniques: learning  
(hidden or latent) structure in  
unlabeled data.

e.g. PCA of legislators's votes: want to see  
how they are organized—by party? by  
ideology? by race?

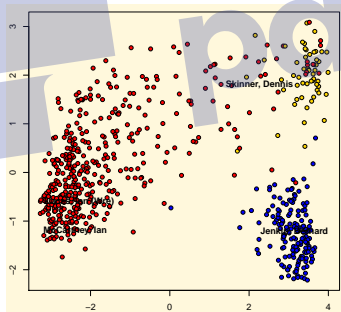


# Remember...

Remove Watermark Now

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



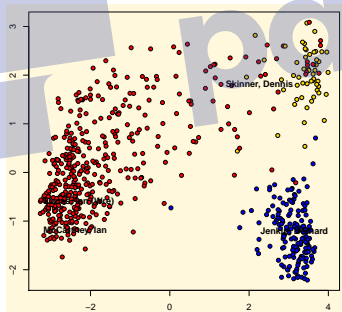
Supervised techniques:

# Remember...

Remove Watermark Now

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



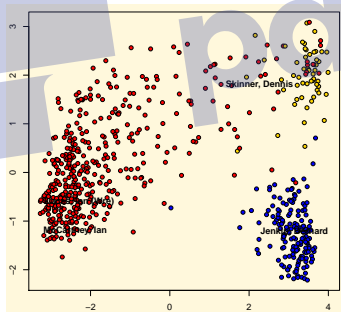
Supervised techniques: learning relationship between inputs and a labeled set of outputs.

# Remember...

Remove Watermark Now

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



Supervised techniques: learning relationship between inputs and a labeled set of outputs.

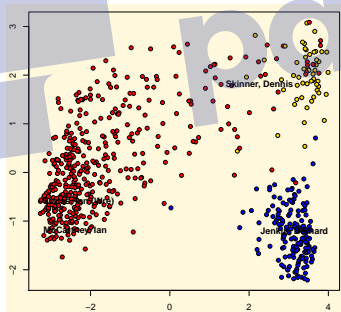
e.g. opinion mining:

# Remember...

Remove Watermark Now

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

- e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



Supervised techniques: learning relationship between inputs and a labeled set of outputs.

- e.g. opinion mining: what makes a critic like or dislike a movie ( $y_i \in \{0, 1\}$ )?

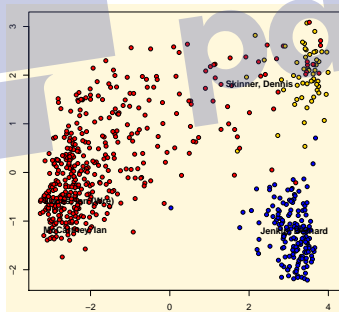


# Remember...

Remove Watermark Now

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

- e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?




Supervised techniques: learning relationship between inputs and a labeled set of outputs.


- e.g. opinion mining: what makes a critic like or dislike a movie ( $y_i \in \{0, 1\}$ )?


**CRITIC REVIEWS FOR STAR WARS: EPISODE VII - THE FORCE AWAKENS**

All Critics (313) | Top Critics (48) | My Critics | Fresh (293) | Rotten (20)


 The new movie, as an act of pure storytelling, streams by with fluency and zip.


[Full Review...](#) | December 21, 2015

 **Anthony Lane**  
New Yorker  
★ Top Critic


 At the end The Force Awakens looks more like a nostalgic film that will work as a transition to the new Star Wars' age. [Full Review in Spanish]


[Full Review...](#) | December 29, 2015

 **Salvador Franco Reyes**

 While Star Wars: The Force Awakens gets temporarily bogged down taking us back to the world that we left in 1983, it introduces us to the new and exciting torch-bearers of the franchise.

[Full Review...](#) | December 30, 2015

 **Blake Howard**  
Graffiti With Punctuation

 This film is a well-planned product that balances nostalgia with the capacity to attract new generations into the Star Wars universe. [Full Review in Spanish]

[Full Review...](#) | December 29, 2015

## Workflow of Supervised Techniques



# Workflow of Supervised Techniques

[Remove Watermark Now](#)

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).

pdfelement

# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**,

# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**, possibly matched in some way to nature of problem.

# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**, possibly matched in some way to nature of problem.
- 3 Adjust algorithm for **optimal performance**,

# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**, possibly matched in some way to nature of problem.
- 3 Adjust algorithm for **optimal performance**, perhaps using **validation set** and/or some kind of **cross-validation**.

# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**, possibly matched in some way to nature of problem.
- 3 Adjust algorithm for **optimal performance**, perhaps using **validation set** and/or some kind of **cross-validation**.
- 4 Report accuracy in **test set**,



# Workflow of Supervised Techniques

Remove Watermark Now

- 1 Obtain/code the **training set** and decide on relevant features (preprocess).
- 2 Decide on the **algorithm**, possibly matched in some way to nature of problem.
- 3 Adjust algorithm for **optimal performance**, perhaps using **validation set** and/or some kind of **cross-validation**.
- 4 Report accuracy in **test set**, possibly combine with other learners in **ensemble**.



Supervised techniques are about **learning** relationship between  $X$  and labeled data.

pdfelement

Supervised techniques are about **learning** relationship between  $X$  and labeled data. Often used interchangeably with **machine learning**: idea that computers could 'learn' relationships without specific programming.



Supervised techniques are about **learning** relationship between  $X$  and labeled data. Often used interchangeably with **machine learning**: idea that computers could 'learn' relationships without specific programming.

Typically used when  $p \gg n$ :

Supervised techniques are about **learning** relationship between  $X$  and labeled data. Often used interchangeably with **machine learning**: idea that computers could 'learn' relationships without specific programming.

Typically used when  $p \gg n$ : number of parameters (e.g. coefficients on terms) is far larger than number of observations (e.g. speakers or documents).

Supervised techniques are about **learning** relationship between  $X$  and labeled data. Often used interchangeably with **machine learning**: idea that computers could 'learn' relationships without specific programming.

Typically used when  $p \gg n$ : number of parameters (e.g. coefficients on terms) is far larger than number of observations (e.g. speakers or documents).

→ results in general **curse of dimensionality** wherein feature matrix is large (e.g. 100k columns) and **sparse** and thus obtaining meaningful estimates is difficult.

Supervised techniques are about **learning** relationship between  $X$  and labeled data. Often used interchangeably with **machine learning**: idea that computers could 'learn' relationships without specific programming.

Typically used when  $p \gg n$ : number of parameters (e.g. coefficients on terms) is far larger than number of observations (e.g. speakers or documents).

→ results in general **curse of dimensionality** wherein feature matrix is large (e.g. 100k columns) and **sparse** and thus obtaining meaningful estimates is difficult.

So techniques may require careful tuning of **regularization parameters** to obtain good performance.





# Notes and Issues II

Once we have the **training set** we face a dilemma...

[Remove Watermark Now](#)

pdfelement

# Notes and Issues II

Once we have the **training set** we face a dilemma...

[Remove Watermark Now](#)

- 1 we can use our technique to fit a very complicated model to this data (perfectly),

pdfelement

# Notes and Issues II

Once we have the **training set** we face a dilemma...

[Remove Watermark Now](#)

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements.

pdfelement

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).



pdfelement

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set,

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).
- we have **overfit** to our training set, and may do poorly on our test set.

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set,



# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance.

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**,

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**, but may induce **bias** in the sense that we miss important relationships in the data.

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**, but may induce **bias** in the sense that we miss important relationships in the data.  
→ we have **underfit** to our training set,

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**, but may induce **bias** in the sense that we miss important relationships in the data.  
→ we have **underfit** to our training set, and may do poorly on our test set.

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
- 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**, but may induce **bias** in the sense that we miss important relationships in the data.  
→ we have **underfit** to our training set, and may do poorly on our test set.

So managing the ***bias-variance tradeoff*** is a key element of supervised learning,

# Notes and Issues II

Remove Watermark Now

Once we have the **training set** we face a dilemma...

- 1 we can use our technique to fit a very complicated model to this data (perfectly), including noisy elements. This avoids **bias**, but it incurs **variance** (when we move to the test set).  
→ we have **overfit** to our training set, and may do poorly on our test set.
  - 2 we can use be more relaxed about the fit of our algorithm in the training set, and accept some imperfections in performance. This avoids high **variance**, but may induce **bias** in the sense that we miss important relationships in the data.  
→ we have **underfit** to our training set, and may do poorly on our test set.
- So managing the **bias-variance tradeoff** is a key element of supervised learning, and we may need to **tune** our algorithms with that in mind.

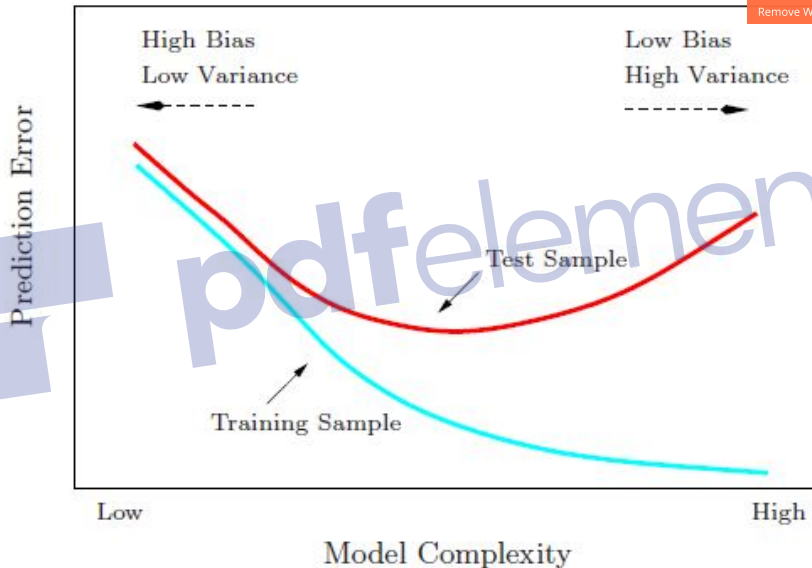
# Bias-Variance Tradeoff (Hastie et al, p38)

[Remove Watermark Now](#)

 pdfelement



# Bias-Variance Tradeoff (Hastie et al, p38)



Remove Watermark Now

# Partner Exercise

[Remove Watermark Now](#)

 pdfelement

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

He suggests you fit various (high dimensional) models to this data, and then pick the one with the lowest prediction error to fit to future data.

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

He suggests you fit various (high dimensional) models to this data, and then pick the one with the lowest prediction error to fit to future data. When you do so, you are unable to obtain anything like the same performance for the arriving 2016 data.

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

He suggests you fit various (high dimensional) models to this data, and then pick the one with the lowest prediction error to fit to future data. When you do so, you are unable to obtain anything like the same performance for the arriving 2016 data. Why?

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

He suggests you fit various (high dimensional) models to this data, and then pick the one with the lowest prediction error to fit to future data. When you do so, you are unable to obtain anything like the same performance for the arriving 2016 data. Why? What's wrong with this approach?

# Partner Exercise

Remove Watermark Now

You are working for `yelp.com`.

Your boss wants to know what will predict reviews for restaurants in Greenwich Village in 2016, and gives you all the data from the Village in 2015.

He suggests you fit various (high dimensional) models to this data, and then pick the one with the lowest prediction error to fit to future data. When you do so, you are unable to obtain anything like the same performance for the arriving 2016 data. Why? What's wrong with this approach?



# Cross-Validation

[Remove Watermark Now](#)

 pdfelement

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

pdfelement

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.



pdfelement

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- 1 divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'),

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- 1 divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- 1 divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- 2 repeat the following  $k$  times (folds):

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

[Remove Watermark Now](#)

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ② grab the other  $k - 1$  chunks as a **training** set



# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

Remove Watermark Now

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ② grab the other  $k - 1$  chunks as a **training set**
  - ③ test on the validation set,

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

Remove Watermark Now

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ② grab the other  $k - 1$  chunks as a **training set**
  - ③ test on the validation set, record **prediction error**

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

Remove Watermark Now

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ② grab the other  $k - 1$  chunks as a **training set**
  - ③ test on the validation set, record **prediction error**
- ③ Average over runs to get prediction error estimate.

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

Remove Watermark Now

Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ① divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ② repeat the following  $k$  times (folds):
  - ① grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ② grab the other  $k - 1$  chunks as a **training set**
  - ③ test on the validation set, record **prediction error**
- ③ Average over runs to get prediction error estimate.

→ Can follow same steps for models of different specifications; variant on this approach can be used for **model selection**, directly.

# Cross-Validation

Want to properly **estimate model performance** (bias and esp. variance).

Remove Watermark Now

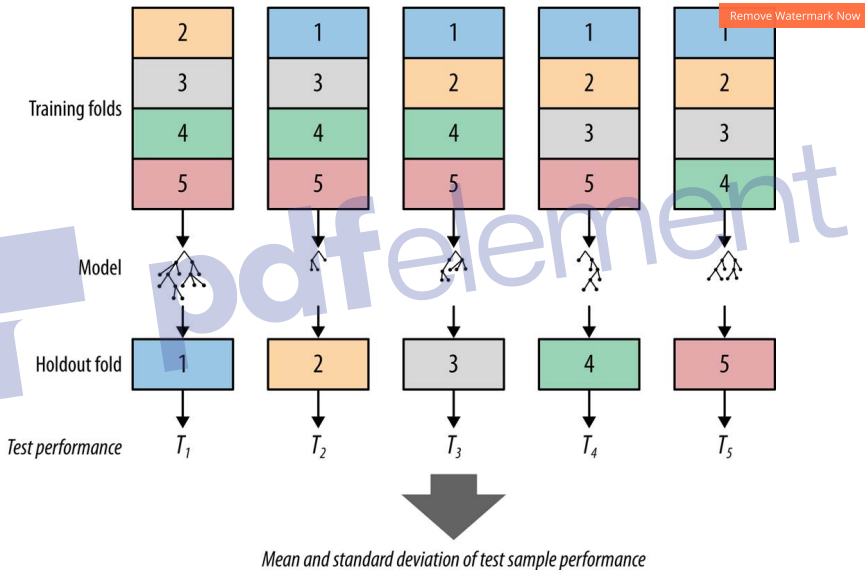
Popular and efficient approach is  **$k$ -fold cross validation**, esp when we don't have enough data to form a completely separate **validation** set.

- ➊ divide all data into  $k$  equal size chunks ( $k = 10$  is common;  $k = n$  is 'leave one out'), and set the parameter(s) of model at particular value,
- ➋ repeat the following  $k$  times (folds):
  - ➊ grab one of the  $k$  chunks as a **validation set** (each only used once)
  - ➋ grab the other  $k - 1$  chunks as a **training set**
  - ➌ test on the validation set, record **prediction error**
- ➌ Average over runs to get prediction error estimate.

→ Can follow same steps for models of different specifications; variant on this approach can be used for **model selection**, directly.



# Graphically



# Motivating Example Diermeier et al, 2011

[Remove Watermark Now](#)

 pdfelement



# Motivating Example Diermeier et al, 2011

Remove Watermark Now



# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.



# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004),



# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004), selected based on their **NOMINATE** scores



# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004), selected based on their **NOMINATE** scores (from roll call behavior, only)



# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004), selected based on their **NOMINATE** scores (from roll call behavior, only)

Code the Republicans as +1, Dems as -1:

# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004), selected based on their **NOMINATE** scores (from roll call behavior, only)

Code the Republicans as +1, Dems as -1:  
 $y_i \in \{-1, +1\}$

# Motivating Example Diermeier et al, 2011

Remove Watermark Now

Diermeier et al want to know what **terms** are most indicative of conservative or liberal positions in legislative debates.

Study the speeches of the 25 **most liberal** and the 25 **most conservative** senators (1989–2004), selected based on their **NOMINATE** scores (from roll call behavior, only)

Code the Republicans as +1, Dems as -1:  
 $y_i \in \{-1, +1\}$

Have the (stemmed, stopped, weighted etc) **speech term matrix** for each Senator as  $X$ .



# Methodological Problem

[Remove Watermark Now](#)

 pdfelement

# Methodological Problem

[Remove Watermark Now](#)

Want to know the relationship between using different terms and ideological views.

pdfelement

# Methodological Problem

[Remove Watermark Now](#)

Want to know the relationship between using different terms and ideological views.

e.g. is 'Ethanol' a 'Democrat' term or a 'Republican' term,

pdfelement

# Methodological Problem

Remove Watermark Now

Want to know the relationship between using different terms and ideological views.

e.g. is 'Ethanol' a 'Democrat' term or a 'Republican' term, and to what degree?



pdfelement

# Methodological Problem

Remove Watermark Now

Want to know the relationship between using different terms and ideological views.

e.g. is 'Ethanol' a 'Democrat' term or a 'Republican' term, and to what degree?

Their **training** set is speech output of most extreme Senators between the 101st and 107th Congress

# Methodological Problem

Remove Watermark Now

Want to know the relationship between using different terms and ideological views.

e.g. is 'Ethanol' a 'Democrat' term or a 'Republican' term, and to what degree?

Their **training** set is speech output of most extreme Senators between the 101st and 107th Congress

Their **test** set is speech output of most extreme Senators in 108th Congress.

# Methodological Problem

Remove Watermark Now

Want to know the relationship between using different terms and ideological views.

e.g. is 'Ethanol' a 'Democrat' term or a 'Republican' term, and to what degree?

Their **training** set is speech output of most extreme Senators between the 101st and 107th Congress

Their **test** set is speech output of most extreme Senators in 108th Congress.

What method to use?

# Support Vector Machines

[Remove Watermark Now](#)

 pdfelement



# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

[Remove Watermark Now](#)

pdfelement

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

[Remove Watermark Now](#)

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.

pdfelement

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches.

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

[Remove Watermark Now](#)

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ :

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

[Remove Watermark Now](#)

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**:

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

Remove Watermark Now

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**: if we plot the Senators in two dimensions  $(x_1, x_2)$ , we can divide them (perfectly) into the two parties using a **straight line**.

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

Remove Watermark Now

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**: if we plot the Senators in two dimensions  $(x_1, x_2)$ , we can divide them (perfectly) into the two parties using a **straight line**.

→ of course,



# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

Remove Watermark Now

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**: if we plot the Senators in two dimensions  $(x_1, x_2)$ , we can divide them (perfectly) into the two parties using a **straight line**.

→ of course, a real problem would have  $p$  being large,

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

Remove Watermark Now

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**: if we plot the Senators in two dimensions  $(x_1, x_2)$ , we can divide them (perfectly) into the two parties using a **straight line**.

→ of course, a real problem would have  $p$  being large, and the space being very high dimensional,

# Support Vector Machines

**Idea** a classifier that builds a model from a binary labeled training set and returns an optimal **hyperplane** that puts new examples into one of the categories non-probabilistically.

Remove Watermark Now

**Here** We have 10 Senators that belong to one of two classes: five Republicans, five Democrats.  $y_i \in \{-1, +1\}$   
Each senator has a number of  $p$  features, which are the term weights from their speeches. To make things simple, suppose that  $p = 2$ : there are (only) two features,  $x_1$  and  $x_2$  per observation.

Assume that the observations are **linearly separable**: if we plot the Senators in two dimensions  $(x_1, x_2)$ , we can divide them (perfectly) into the two parties using a **straight line**.

→ of course, a real problem would have  $p$  being large, and the space being very high dimensional, but the logic is the same.

# The 10 Senators

[Remove Watermark Now](#)

 pdfelement

# The 10 Senators

Remove Watermark Now



# The 10 Senators

Remove Watermark Now



As the parties linearly separably?

# The 10 Senators

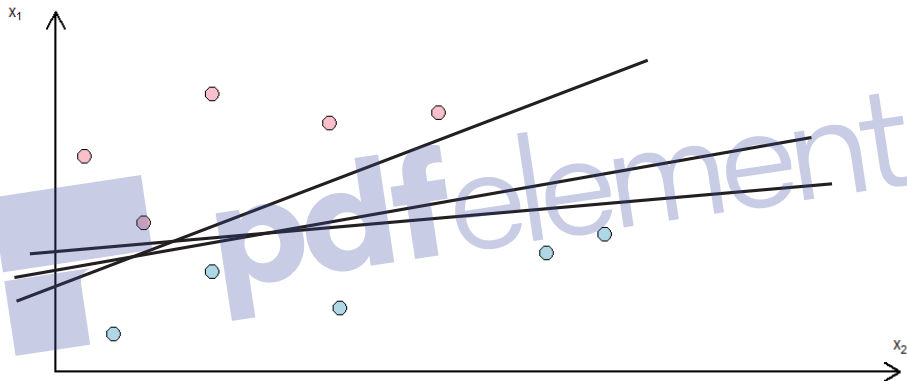
Remove Watermark Now



As the parties linearly separably? Where could you draw the line?

# The 10 Senators

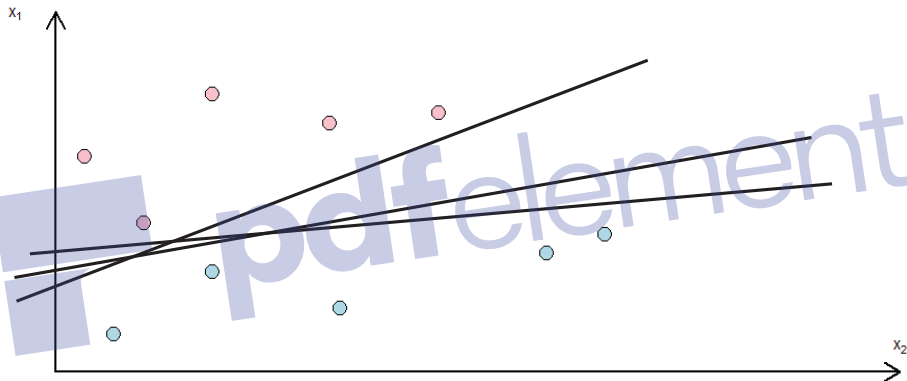
Remove Watermark Now





# The 10 Senators

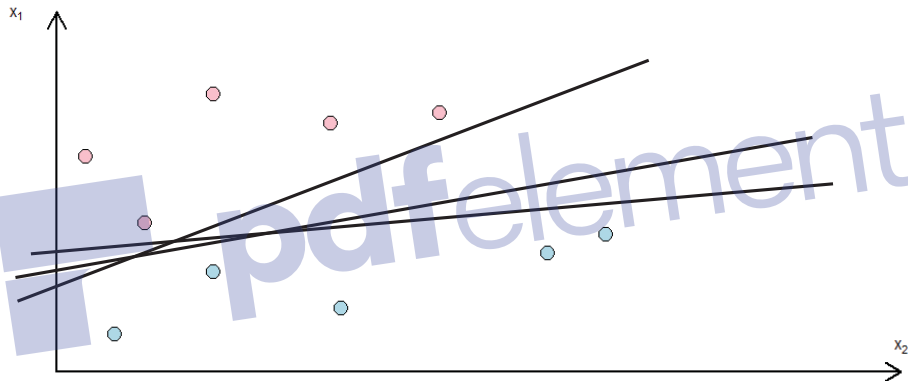
Remove Watermark Now



Which line should we prefer?

# The 10 Senators

Remove Watermark Now



Which line should we prefer?

# What is the Optimal Hyperplane?

Remove Watermark Now

 pdfelement

# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

pdfelement

# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

But we want to avoid ones that pass close to the points:

pdfelement

# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

But we want to avoid ones that pass close to the points: such lines will tend to be sensitive to noise and make classification errors with future examples.

# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

But we want to avoid ones that pass close to the points: such lines will tend to be sensitive to noise and make classification errors with future examples.

So pick line that gives largest minimum distance from the training cases.

# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

But we want to avoid ones that pass close to the points: such lines will tend to be sensitive to noise and make classification errors with future examples.

So pick line that gives largest minimum distance from the training cases. That is, the line that's as far as possible from the closest cases on both sides.



# What is the Optimal Hyperplane?

Remove Watermark Now

There are many possible lines.

But we want to avoid ones that pass close to the points: such lines will tend to be sensitive to noise and make classification errors with future examples.

So pick line that gives largest minimum distance from the training cases. That is, the line that's as far as possible from the closest cases on both sides.

→ That optimal line—the separating hyperplane—is the maximum margin hyperplane. It will maximize the margin of the training data.

# Partner Exercise

[Remove Watermark Now](#)

 pdfelement

# Partner Exercise

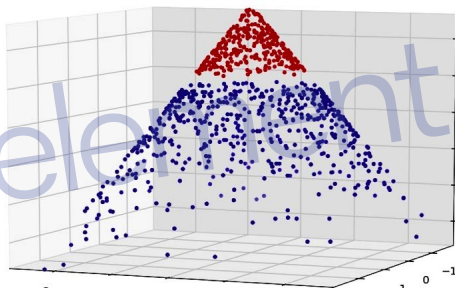
Remove Watermark Now

Consider the figure.

It's a situation where each Senator's features are of three dimensions (rather than two).

How could we (optimally) separate the data in a linear way?

Can we still use a line?



from <http://www.edvancer.in/>

# How to get it I: Notation Variants

[Remove Watermark Now](#)

 pdfelement

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

pdfelement

# How to get it I: Notation Variants

[Remove Watermark Now](#)

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning,

pdfelement

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

pdfelement

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

$$\rightarrow \mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y$$



# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

$$\rightarrow \mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0.$$

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0.$

or use dot product:

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular:

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular: can use a two dimensional form—

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular: can use a two dimensional form—

$\mathbf{w} = [-a, 1]$  and  $\mathbf{x} = [x, y]$ ,

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular: can use a two dimensional form—

$\mathbf{w} = [-a, 1]$  and  $\mathbf{x} = [x, y]$ , so that  $\mathbf{w} \cdot \mathbf{x} = y - ax$

# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular: can use a two dimensional form—

$\mathbf{w} = [-a, 1]$  and  $\mathbf{x} = [x, y]$ , so that  $\mathbf{w} \cdot \mathbf{x} = y - ax$

thus  $\mathbf{w} \cdot \mathbf{x} - b = 0$



# How to get it I: Notation Variants

Remove Watermark Now

We can write any line as  $y = ax + b$  or  $y - ax - b = 0$ .

In machine learning, often rewrite terms as  $\mathbf{w} = [-b, -a, 1]$  and  $\mathbf{x} = [1, x, y]$

→  $\mathbf{w}^T \mathbf{x} = -b \times 1 + (-a) \times x + 1 \times y = y - ax - b = 0$ .

or use dot product:  $\mathbf{w} \cdot \mathbf{x} = 0$

Also popular: can use a two dimensional form—

$\mathbf{w} = [-a, 1]$  and  $\mathbf{x} = [x, y]$ , so that  $\mathbf{w} \cdot \mathbf{x} = y - ax$

thus  $\mathbf{w} \cdot \mathbf{x} - b = 0$  (since  $0 = y - ax - b$ )

# How to get it II: working with the equation

[Remove Watermark Now](#)

 pdfelement

# How to get it II: working with the equation

[Remove Watermark Now](#)

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

pdfelement

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is equidistant,

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is **equidistant**, i.e. half-way between,

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

Those parallel lines should be as far from each other as possible without misclassifying anybody:

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

Those parallel lines should be as far from each other as possible without misclassifying anybody: the distance between them is the **margin** and we want that to be **maximized**.



# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy

$$\mathbf{w} \cdot \mathbf{x} - b = 0$$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

Those parallel lines should be as far from each other as possible without misclassifying anybody: the distance between them is the **margin** and we want that to be **maximized**.

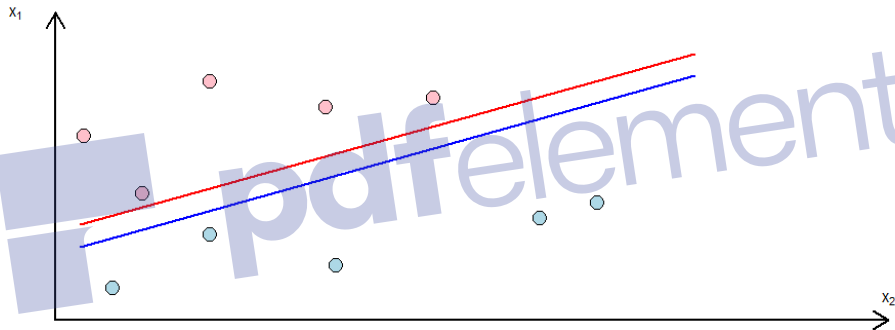
# How Could We Do Better?

[Remove Watermark Now](#)

 pdfelement

# How Could We Do Better?

Remove Watermark Now



# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy  
 $\mathbf{w} \cdot \mathbf{x} - b = 0$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

Those parallel lines should be as far from each other as possible without misclassifying anybody: the distance between them is the **margin** and we want that to be **maximized**.

# How to get it II: working with the equation

Remove Watermark Now

Our hyperplane (line) will separate the data, and will satisfy  
 $\mathbf{w} \cdot \mathbf{x} - b = 0$

We can think of it as the line that is **equidistant**, i.e. half-way between, two **parallel** hyperplanes ( $H_R$  and  $H_D$ ) that separate the Reps from the Dems.

Those parallel lines should be as far from each other as possible without misclassifying anybody: the distance between them is the **margin** and we want that to be **maximized**.

**NB** The hyperplane cannot be anywhere other than **equidistant** because then it will break the rule about ensuring the **largest minimum distance**.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

# The parallel hyperplanes

[Remove Watermark Now](#)

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $x$ s).

pdfelement

# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $x$ s).

This means that 'their' hyperplane should 'capture' them all, and separate them fully from the Democrats (in our 2D space).



# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $x$ s).

This means that 'their' hyperplane should 'capture' them all, and separate them fully from the Democrats (in our 2D space).

→ requires that  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = 1$

# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $x$ s).

This means that 'their' hyperplane should 'capture' them all, and separate them fully from the Democrats (in our 2D space).

→ requires that  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = 1$

Same idea for the Democrats:

# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $\mathbf{x}$ s).

This means that 'their' hyperplane should 'capture' them all, and separate them fully from the Democrats (in our 2D space).

→ requires that  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = 1$

Same idea for the Democrats:  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = -1$

# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $\mathbf{x}$ s).

This means that 'their' hyperplane should 'capture' them all, and separate them fully from the Democrats (in our 2D space).

→ requires that  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = 1$

Same idea for the Democrats:  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = -1$

The distance between the two hyperplanes ( $H_R$  and  $H_D$ ) we construct is the **margin**,

# The parallel hyperplanes

Remove Watermark Now

We want all the Republicans ( $y_i = 1$ )—and *only* the Republicans—to be classified as Republicans (given their  $\mathbf{x}$ s).

This means that ‘their’ hyperplane should ‘capture’ them all, and separate them fully from the Democrats (in our 2D space).

→ requires that  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = 1$

Same idea for the Democrats:  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = -1$

The distance between the two hyperplanes ( $H_R$  and  $H_D$ ) we construct is the **margin**, and its width is  $\frac{2}{\|\mathbf{w}\|}$ , where  $\|\mathbf{w}\|$  is the norm of  $\mathbf{w}$ .

# Optimization Problem

[Remove Watermark Now](#)

 pdfelement

# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the **constraints** we mentioned. This will ensure the **widest possible** margin.

pdfelement

# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the **constraints** we mentioned. This will ensure the **widest possible** margin.

Well rearranging:

pdfelement



# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the constraints we mentioned. This will ensure the widest possible margin.

Well rearranging: minimize  $\|\mathbf{w}\|$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x} - b) \geq 1 \quad \forall i$

pdfelement

# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the constraints we mentioned. This will ensure the widest possible margin.

Well rearranging: minimize  $\|\mathbf{w}\|$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x} - b) \geq 1 \quad \forall i$

This will give us the two hyperplanes  $H_R$  and  $H_D$ , and the line equidistant between them will be our classifier.

# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the constraints we mentioned. This will ensure the widest possible margin.

Well rearranging: minimize  $\|\mathbf{w}\|$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x} - b) \geq 1 \quad \forall i$

This will give us the two hyperplanes  $H_R$  and  $H_D$ , and the line equidistant between them will be our classifier.

We get the relevant value of  $b$  by noting that  $\frac{b}{\|\mathbf{w}\|}$  is the (perpendicular) distance from the optimal hyperplane to the origin.

# Optimization Problem

Remove Watermark Now

Need to find minimum value of  $\|\mathbf{w}\|$ , subject to the **constraints** we mentioned. This will ensure the **widest possible** margin.

**Well** rearranging: minimize  $\|\mathbf{w}\|$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x} - b) \geq 1 \quad \forall i$

This will give us the two hyperplanes  $H_R$  and  $H_D$ , and the line equidistant between them will be our classifier.

We get the relevant value of  $b$  by noting that  $\frac{b}{\|\mathbf{w}\|}$  is the (perpendicular) distance from the optimal hyperplane to the origin.

Turns out that the minimization of  $\|\mathbf{w}\|$  is amenable to quadratic programming methods.

# Graphically...

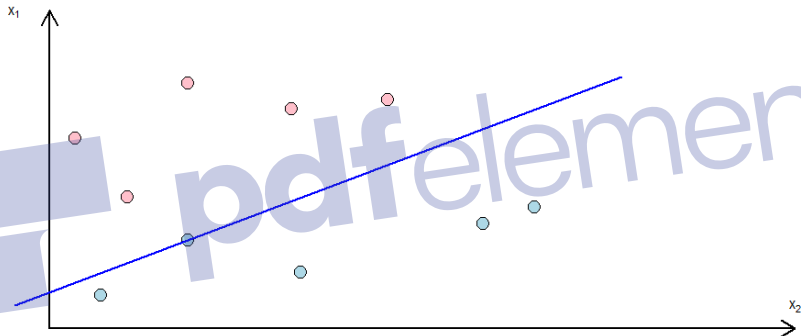
[Remove Watermark Now](#)



pdfelement

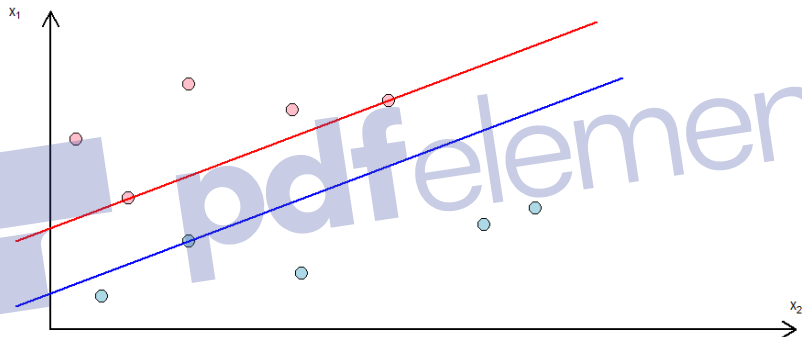
# Graphically...

Remove Watermark Now



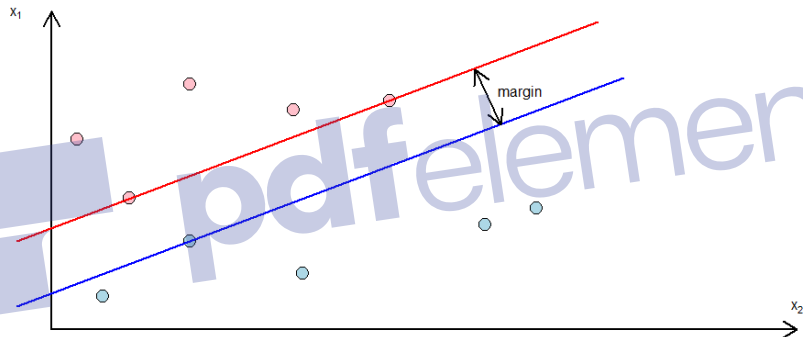
# Graphically...

Remove Watermark Now



# Graphically...

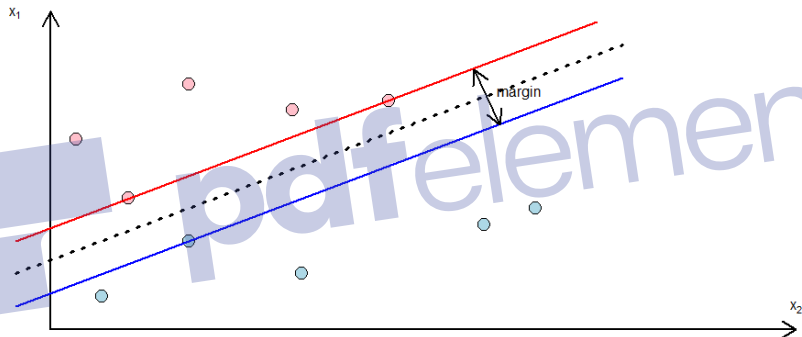
Remove Watermark Now





# Graphically...

[Remove Watermark Now](#)





# Notes

Remove Watermark Now

At least one of the Senators, in one of the parties, will be on their line.

pdfelement

At least one of the Senators, in one of the parties, will be on their line. Why?

pdfelement

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line.

pdfelement

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger....

# Notes

Remove Watermark Now

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines,

# Notes

Remove Watermark Now

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines, we've broken the rule about no misclassifications,



# Notes

Remove Watermark Now

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines, we've broken the rule about no misclassifications, and that's a constraint we have to follow.

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines, we've broken the rule about no misclassifications, and that's a constraint we have to follow.

In fact, the points closest to the separating hyperplane will be the Senators lying on their (respective) parallel hyperplanes.

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines, we've broken the rule about no misclassifications, and that's a constraint we have to follow.

In fact, the points closest to the separating hyperplane will be the Senators lying on their (respective) parallel hyperplanes.

We use the term support vectors to describe the training examples closest to the hyperplane.

At least one of the Senators, in one of the parties, will be on their line. Why?

A: Suppose nobody is on either line. Then we must be able to move the lines 'apart' a little further and make the margin bigger...

But as soon as someone is pushed between the lines, we've broken the rule about no misclassifications, and that's a constraint we have to follow.

In fact, the points closest to the separating hyperplane will be the Senators lying on their (respective) parallel hyperplanes.

We use the term support vectors to describe the training examples closest to the hyperplane. Those support vectors completely determine where our maximum-margin hyperplane will be.

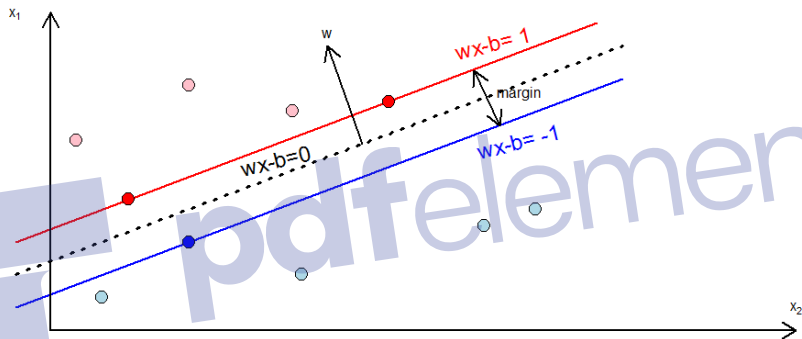
# The Support Vectors

[Remove Watermark Now](#)

 pdfelement

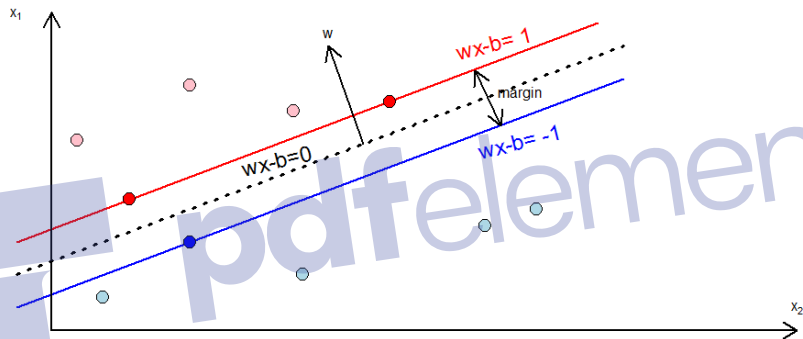
# The Support Vectors

Remove Watermark Now



# The Support Vectors

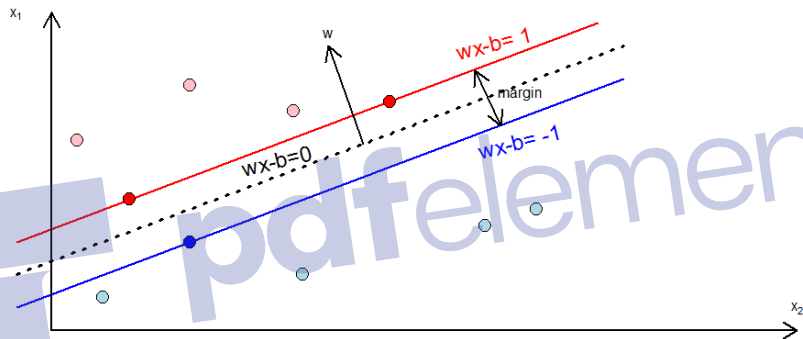
Remove Watermark Now



The support vectors lie on...

# The Support Vectors

Remove Watermark Now



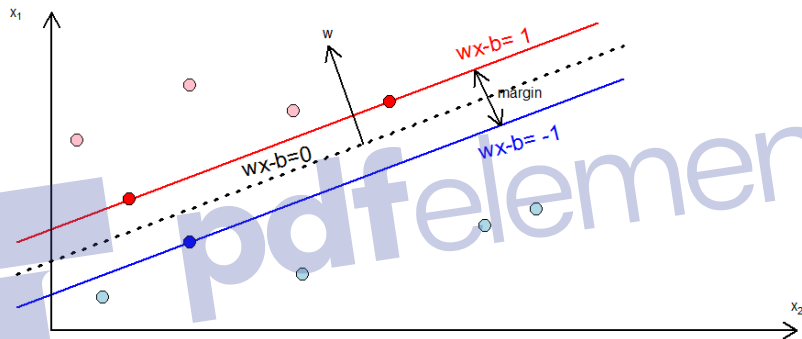
The support vectors lie on...

$$\mathbf{w} \cdot \mathbf{x} - b = 1$$



# The Support Vectors

Remove Watermark Now



The support vectors lie on...

$$\mathbf{w} \cdot \mathbf{x} - b = 1 \text{ or } \mathbf{w} \cdot \mathbf{x} - b = -1$$

# SVM weights

[Remove Watermark Now](#)

 pdfelement

# SVM weights

So We have a hyperplane that separates the Democrats from the Republicans.

[Remove Watermark Now](#)

pdfelement

# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.



pdfelement

# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.

i.e. if the product is positive  $\rightarrow$  **Republican**; if negative  $\rightarrow$  **Democrat**

# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.

i.e. if the product is positive  $\rightarrow$  **Republican**; if negative  $\rightarrow$  **Democrat**

Plus for each feature,  $x_1, x_2, \dots$ , the vector  $\mathbf{w}$  gives us a **weight**.

# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.

i.e. if the product is positive  $\rightarrow$  **Republican**; if negative  $\rightarrow$  **Democrat**

Plus for each feature,  $x_1, x_2, \dots$ , the vector  $\mathbf{w}$  gives us a **weight**. The absolute size of the weights—relative to one another—is a measure of **how important** that feature is for separating the Senators into Democrat and Republican.

# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.

i.e. if the product is positive  $\rightarrow$  **Republican**; if negative  $\rightarrow$  **Democrat**

Plus for each feature,  $x_1, x_2, \dots$ , the vector  $\mathbf{w}$  gives us a **weight**. The absolute size of the weights—relative to one another—is a measure of **how important** that feature is for separating the Senators into Democrat and Republican.

NB SVMs typically have few features with **non-zero** weights,



# SVM weights

Remove Watermark Now

So We have a hyperplane that separates the Democrats from the Republicans.

The vector  $\mathbf{w}$  is orthogonal to that, and when we multiply it (dot product) by  $\mathbf{x}$  (and add  $b$  term) we get the **predicted class** of a 'new' Senator.

i.e. if the product is positive  $\rightarrow$  **Republican**; if negative  $\rightarrow$  **Democrat**

Plus for each feature,  $x_1, x_2, \dots$ , the vector  $\mathbf{w}$  gives us a **weight**. The absolute size of the weights—relative to one another—is a measure of **how important** that feature is for separating the Senators into Democrat and Republican.

NB SVMs typically have few features with **non-zero** weights, and those that are non-zero come from the **support vectors** (the 'important' observations)



# Back to Diermeier et al, 2011

Achieve 92% accuracy (!)

[Remove Watermark Now](#)

 pdfelement

# Back to Diermeier et al, 2011

[Remove Watermark Now](#)

Achieve 92% accuracy (!)

Sort words according to coefficients: very positive weights imply **conservative** words; very negative weights imply **liberal** words

pdfelement

# Back to Diermeier et al, 2011

[Remove Watermark Now](#)

Achieve 92% accuracy (!)

Sort words according to coefficients: very positive weights imply **conservative** words; very negative weights imply **liberal** words. Argue that it is 'values' rather than economics that separates liberals from conservatives.

pdfelement

# Back to Diermeier et al, 2011

Remove Watermark Now

Achieve 92% accuracy (!)

Sort words according to coefficients: very positive weights imply **conservative** words; very negative weights imply **liberal** words. Argue that it is 'values' rather than economics that separates liberals from conservatives.

Words			
Liberal		Conservative	
FAS: -199.49	SBA: -113.10	habeas: 193.55	homosexual: 103.07
Ethanol: -198.92	Nursing: -109.38	CFTC: 187.16	everglades: 102.87
Wealthiest: -159.74	Providence: -108.73	surtax: 151.81	tower: 101.67
Collider: -142.28	Arctic: -108.30	marriage: 145.79	tripartisan: 101.23
WIC: -140.14	Orange: -107.98	cloning: 141.71	PRC: 102.90
ILO: -139.89	Glaxo: -107.81	tritium: 133.49	scouts: 97.55
Handgun: -129.01	Libraries: -107.70	ranchers: 132.95	nashua: 99.32
Lobbyists: -128.95	Disabilities: -106.44	BTU: 121.92	ballistic: 97.22
Enron: -127.71	Prescription: -106.31	grazing: 121.59	salting: 94.28
Fishery: -127.30	NIH: -105.52	unfunded: 120.82	abortion: 91.94
Hydrogen: -122.59	Lobbying: -105.35	catfish: 120.82	NTSB: 93.81
Souter: -121.40	NRA: -105.20	IRS: 114.91	Haiti: 97.28
PTSD: -119.87	Trident: -104.15	unborn: 111.88	PAC: 92.85
Gun: -119.52	RNC: -103.46	Taiwan: 111.13	taxing: 90.39

# Partner Exercise

[Remove Watermark Now](#)

 pdfelement

# Partner Exercise

[Remove Watermark Now](#)

Diermeier et al's results imply that liberals use 'Handgun' more often than conservatives.





# Partner Exercise

Remove Watermark Now

Diermeier et al's results imply that liberals use 'Handgun' more often than conservatives.

- 1 Does that imply that making conservative Senators use the word 'handgun' more often will make them more liberal?

Diermeier et al's results imply that liberals use 'Handgun' more often than conservatives.

- 1 Does that imply that making conservative Senators use the word 'handgun' more often will make them more liberal? What does your answer suggest about [prediction](#) vs [explanation](#) with supervised techniques?

# Partner Exercise

Remove Watermark Now

Diermeier et al's results imply that liberals use 'Handgun' more often than conservatives.

- 1 Does that imply that making conservative Senators use the word 'handgun' more often will make them more liberal? What does your answer suggest about **prediction** vs **explanation** with supervised techniques?
- 2 what is the (most likely) problem in the causal claim that  $X \rightarrow Y$  in the Diermeier et al study?

Wait...

[Remove Watermark Now](#)

 pdfelement

# Wait...

[Remove Watermark Now](#)

What's so special about SVM?

pdfelement

# Wait...

[Remove Watermark Now](#)

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression,

pdfelement

# Wait...

[Remove Watermark Now](#)

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression, regularized in some way (RLR) to avoid over-fitting?

pdfelement

# Wait...

[Remove Watermark Now](#)

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression, regularized in some way (RLR) to avoid over-fitting?

→ yes, this can work well and is commonly used in machine learning.



# Wait...

[Remove Watermark Now](#)

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression, regularized in some way (RLR) to avoid over-fitting?

→ yes, this can work well and is commonly used in machine learning.

But if you want to know the decision boundary then SVM might be a better choice.

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression, regularized in some way (RLR) to avoid over-fitting?

→ yes, this can work well and is commonly used in machine learning.

But if you want to know the decision boundary then SVM might be a better choice.

→ RLR will optimize probabilities of class membership, rather than just trying to learn the boundary (simpler, more direct task).

# Wait...

Remove Watermark Now

What's so special about SVM? If we have  $y_i \in \{0, 1\}$  why not just use logistic regression, regularized in some way (RLR) to avoid over-fitting?

→ yes, this can work well and is commonly used in machine learning.

But if you want to know the decision boundary then SVM might be a better choice.

→ RLR will optimize probabilities of class membership, rather than just trying to learn the boundary (simpler, more direct task).

BTW RLR can cope well with noise, and (hard margin) SVM will struggle if there is no linear separability...

# What if...

[Remove Watermark Now](#)



pdfelement

# What if...

[Remove Watermark Now](#)

The Senators were **not** linearly separable? ('soft margin' SVM problem)

pdfelement

# What if...

[Remove Watermark Now](#)

The Senators were **not** linearly separable? ('soft margin' SVM problem)

pdfelement

# What if...

[Remove Watermark Now](#)

The Senators were **not** linearly separable? ('soft margin' SVM problem)

pdfelement

# What if...

[Remove Watermark Now](#)

The Senators were **not** linearly separable? ('soft margin' SVM problem)

pdfelement



# Graphically...

Remove Watermark Now



# What if...

[Remove Watermark Now](#)



pdfelement

# What if...

[Remove Watermark Now](#)

The Senators were **not** linearly separable? ('soft margin' SVM problem)

pdfelement

# What if...

Remove Watermark Now

The Senators were **not** linearly separable? ('soft margin' SVM problem)

Can introduce a **hinge loss** function into the minimization problem...

$$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$$

# What if...

Remove Watermark Now

The Senators were **not** linearly separable? ('soft margin' SVM problem)

Can introduce a **hinge loss** function into the minimization problem...

$$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$$

= 0 if the  $x$ s are on the 'correct' side of the margin...

# What if...

Remove Watermark Now

The Senators were **not** linearly separable? ('soft margin' SVM problem)

Can introduce a **hinge loss** function into the minimization problem...

$$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$$

= 0 if the  $x$ s are on the 'correct' side of the margin...

And proportional to the distance from the margin *if* the point is on the 'wrong' side of the margin.

# What if...

Remove Watermark Now

The Senators were **not** linearly separable? ('soft margin' SVM problem)

Can introduce a **hinge loss** function into the minimization problem...

$$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$$

= 0 if the  $x$ s are on the 'correct' side of the margin...

And proportional to the distance from the margin *if* the point is on the 'wrong' side of the margin.

Hyperplane(s) will be drawn in way that is more sensitive to 'bigger' mistakes in classification.

# Partner Exercise

[Remove Watermark Now](#)

 pdfelement



# Partner Exercise

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$

[Remove Watermark Now](#)

pdfelement

## Partner Exercise

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

[Remove Watermark Now](#)

pdfelement

## Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want



pdfelement

# Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

# Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

1 If a classification is correct,

## Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

## Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?
- 2 If a Republican is misclassified as a Democrat, what value will  $\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

## Partner Exercise

Remove Watermark Now

We want  $\mathbf{w} \cdot \mathbf{x} - b \geq 1$ , if  $y_i = \text{Republican}$   
and  $\mathbf{w} \cdot \mathbf{x} - b \leq -1$ , if  $y_i = \text{Democrat}$

Another way (slightly different notation) to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?
- 2 If a Republican is misclassified as a Democrat, what value will  $\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?
- 3 Two Republicans were misclassified as Democrats by the machine. In the first case,  $f(x_i) = -2$ . In the second case,  $f(x_i) = -100$ . Which has the 'worse' value of hinge loss?





Another way to express this is that we want

 pdfelement

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

pdfelement

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

1 If a classification is **correct**,

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?  
Positive,  $\geq 1$ .
- 2 If a Republican is misclassified as a Democrat, what value will  $\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

- 2 If a Republican is misclassified as a Democrat, what value will  $\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Now,  $f(x)$  is negative ('looks like' a Dem), but  $y$  is positive,



Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

- 2 If a Republican is misclassified as a Democrat, what value will

$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Now,  $f(x)$  is negative ('looks like' a Dem), but  $y$  is positive, so  $f(x)y$  is negative.

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

- 2 If a Republican is misclassified as a Democrat, what value will

$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Now,  $f(x)$  is negative ('looks like' a Dem), but  $y$  is positive, so  $f(x)y$  is negative. But then  $1 - f(x)y$  is large, and is the maximum of the set.

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

- 2 If a Republican is misclassified as a Democrat, what value will

$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Now,  $f(x)$  is negative ('looks like' a Dem), but  $y$  is positive, so  $f(x)y$  is negative. But then  $1 - f(x)y$  is large, and is the maximum of the set.

- 3 Two Republicans were misclassified as Democrats by the machine. In the first case,  $f(x_i) = -2$ . In the second case,  $f(x_i) = -100$ . Which has the 'worse' value of hinge loss?

Another way to express this is that we want

$$f(x_i) = \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- 1 If a classification is **correct**, what do we know about the value of  $y_i \times f(x_i)$ ?

Positive,  $\geq 1$ .

- 2 If a Republican is misclassified as a Democrat, what value will

$\mathbb{L}(f(x), y) = \max(0, 1 - f(x)y)$  take?

Now,  $f(x)$  is negative ('looks like' a Dem), but  $y$  is positive, so  $f(x)y$  is negative. But then  $1 - f(x)y$  is large, and is the maximum of the set.

- 3 Two Republicans were misclassified as Democrats by the machine. In the first case,  $f(x_i) = -2$ . In the second case,  $f(x_i) = -100$ . Which has the 'worse' value of hinge loss?

$1 - f(x)y$  is  $1 - (-2)(+1)$  in first case and  $1 - (-100)(+1)$  in second case. Hinge loss larger in second case!

# What if...

[Remove Watermark Now](#)



pdfelement

## What if...

The situation was considerably 'worse',

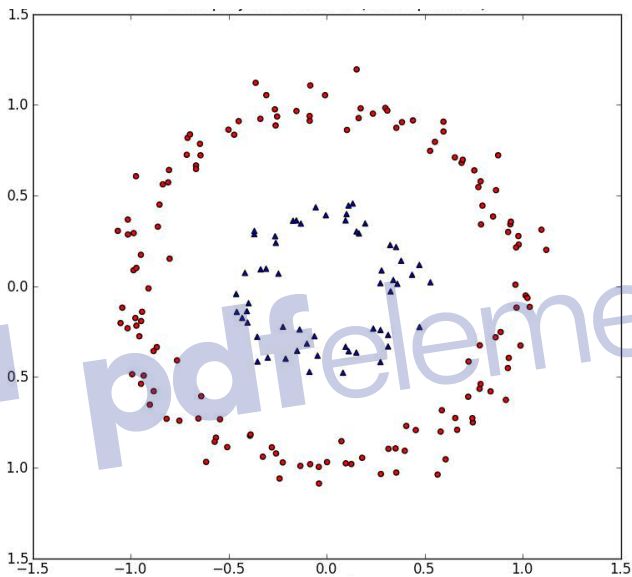


# What if...

[Remove Watermark Now](#)

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

pdfelement



Remove Watermark Now

from [www.eric-kim.net](http://www.eric-kim.net)



# What if...

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

[Remove Watermark Now](#)

pdfelement

# What if...

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

[Remove Watermark Now](#)

May be a way to **transform** the data,



# What if...

[Remove Watermark Now](#)

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ ,

pdfelement

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space,

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space, we can lift it into a **three** dimensional feature space using

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space, we can lift it into a **three** dimensional feature space using

$$\phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2).$$

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space, we can lift it into a **three** dimensional feature space using

$$\phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2).$$

Then use a **linear** SVM on the transformed data set,



# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space, we can lift it into a **three** dimensional feature space using

$$\phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2).$$

**Then** use a **linear** SVM on the transformed data set, and then **map back** to the original 2D space.

# What if...

Remove Watermark Now

The situation was considerably 'worse', such that neither a 'hard margin' nor 'soft margin' linear approach would work?

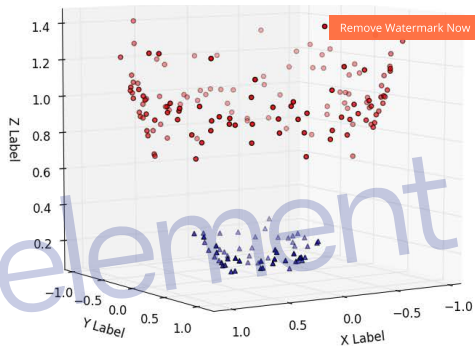
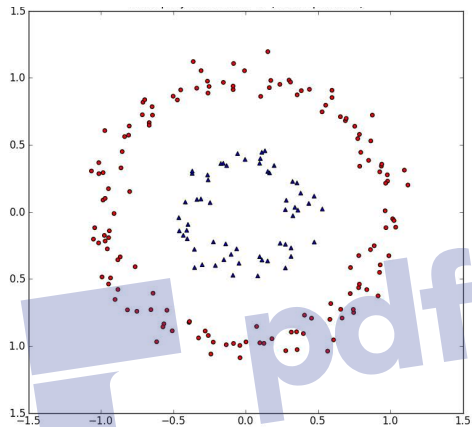
May be a way to **transform** the data, using a transformation  $\phi$ , such that it *can* now be separated using a **linear** classifier.

When the data is in a two dimensional feature space, we can lift it into a **three** dimensional feature space using

$$\phi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2).$$

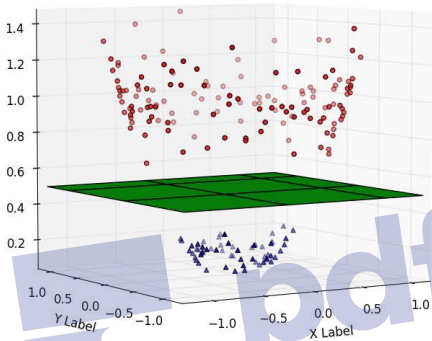
Then use a **linear** SVM on the transformed data set, and then **map back** to the original 2D space.

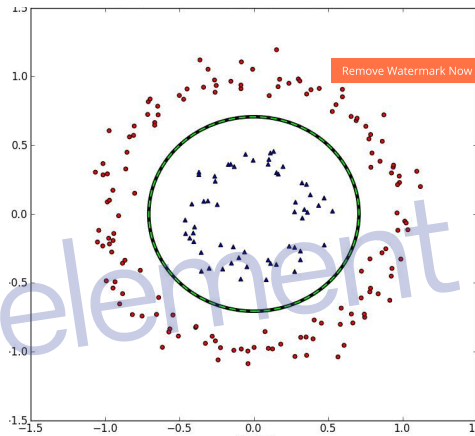
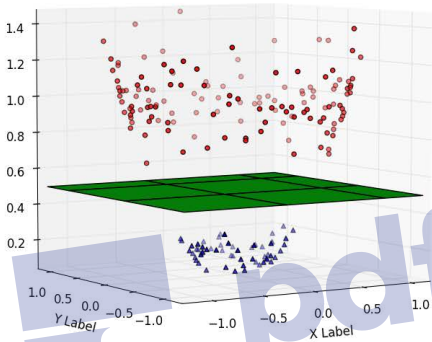
→ Results in a **non-linear** hyperplane once back in 2 dimensions.



from [www.eric-kim.net](http://www.eric-kim.net)

Remove Watermark Now





from [www.eric-kim.net](http://www.eric-kim.net)



# Kernel Methods

Explicitly transforming data into a new space can be very expensive in terms of computation.

[Remove Watermark Now](#)

pdfelement

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?



# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?  
→ exactly what **kernel methods** do:

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed.

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed. Can then feed those inner products to the classifier.

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed. Can then feed those inner products to the classifier.

this '**kernel trick**' cuts cost considerably,

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

- But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?
- exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed. Can then feed those inner products to the classifier.
- this **'kernel trick'** cuts cost considerably, though choosing and tuning the 'correct' kernel may be difficult.

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed. Can then feed those inner products to the classifier.

this 'kernel trick' cuts cost considerably, though choosing and tuning the 'correct' kernel may be difficult.

For text analysis,

# Kernel Methods

Remove Watermark Now

Explicitly transforming data into a new space can be very expensive in terms of computation.

But What if we could do the transformation, and take the distances between observations **implicitly**, and use **them** for our classification?

→ exactly what **kernel methods** do: use **kernel functions**,  $K(\mathbf{x}_i, \mathbf{x}_j)$  to compute the  $i, j$  pairwise **dot products** for training data *as if* it had been transformed. Can then feed those inner products to the classifier.

this 'kernel trick' cuts cost considerably, though choosing and tuning the 'correct' kernel may be difficult.

For text analysis, **string kernels** use a function  $K(a, b)$  to implicitly calculate the distance between strings of characters via the number of subsequences they have in common.



# Partner Exercise

[Remove Watermark Now](#)

 pdfelement

# Partner Exercise

[Remove Watermark Now](#)

Using the ideas we discussed at the start of lecture, how should one go about picking a kernel (from the large variety on offer) for the problem at hand?