

# The Perceptron Algorithm

CSE 446: Machine Learning  
Slides by Emily Fox (mostly)  
Presented by Anna Karlin

April 29, 2019

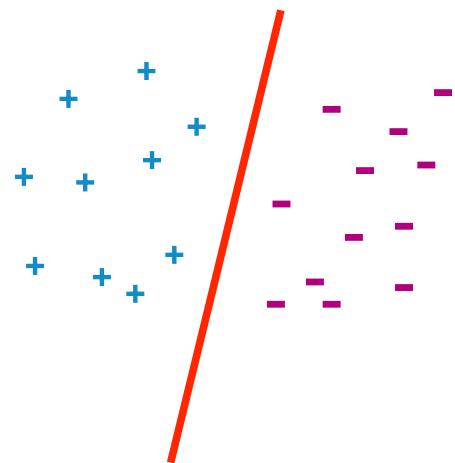
©2017 Emily Fox

# Can we do binary classification without a model?

Before we used the logistic model:

$$P(Y = y | x, w) = 1 / (1 + \exp(-yw^T x))$$

i.e. separate the data into two data groups.  
Find a hyperplane that divides the two.



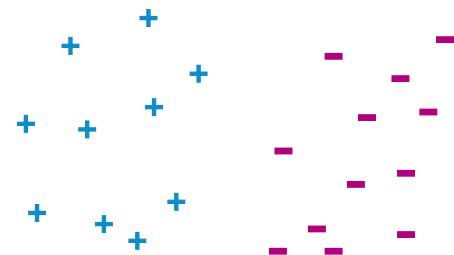
# The perceptron algorithm

# The perceptron algorithm

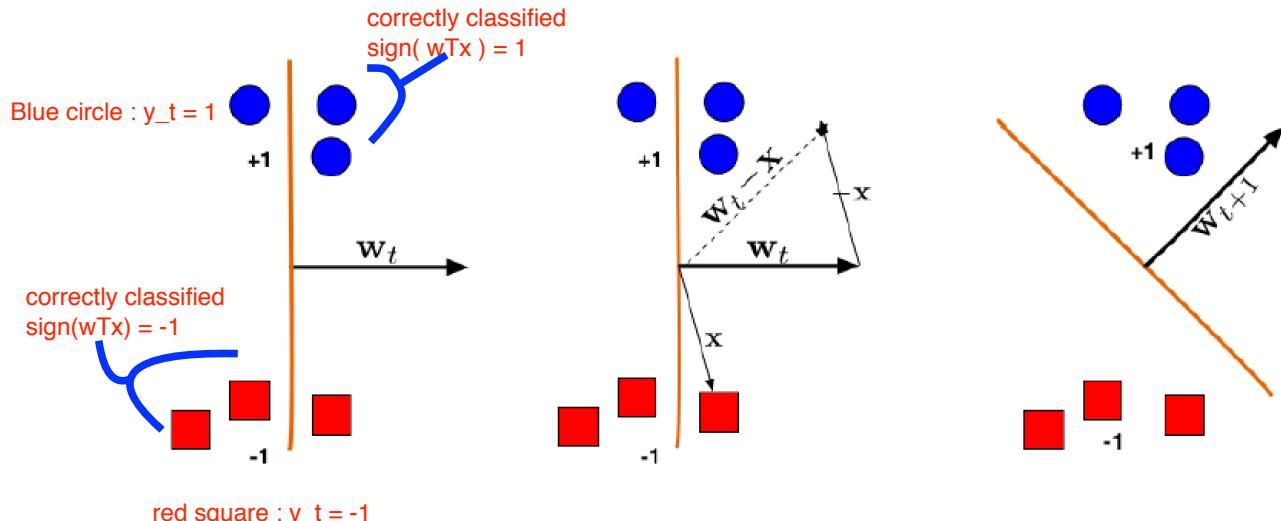
For simplicity in this lecture we will assume that  $b=0$

[Rosenblatt '58, '62]

- Classification setting:  $y$  in  $\{-1, +1\}$
- Linear model
  - Prediction:  
using basic thresholding —> Find a  $w$  and  $b$  based on training set. if  $wTx + b > 0$  choose  $y = +1$ , otherwise choose  $-1$   
 $y_{\text{guess}} = \text{sign}(wTx + b)$
- Training:
  - Initialize weight vector:
  - At each time step:
    - Observe features: find the  $x_t$  of a datapoint in training set (features)
    - Make prediction:  $y_{\text{guess}} = \text{sign}(wTx + b)$
    - Observe true class:  $y_t$  in  $\{-1, 1\}$
    - Update model: ONLINE LEARNING (but could be used in batch setting just one data point at a time)
      - If prediction is not equal to truth  $w_{t+1} = w_t + y_t * x_t$



Picture due to Killian Weinberger



- Initialize  $w_0$ .
  - for  $t := 1$  to  $T$ 
    - Observe feature vector  $x_t$ .
    - Make a prediction:  $\hat{y} = \text{sign}(x_t^T w_t)$
    - Observe the true label  $y_t \in \{-1, 1\}$ .
    - Update the model:
      - \* If  $\hat{y} = y_t$ , then  $w_{t+1} := w_t$ .
      - \* Otherwise ( $\hat{y} \neq y_t$ ), then  $w_{t+1} := w_t + y_t x_t$ .
- if we are incorrect:  
 $w_t + x_t$  if  $y_t = 1$   
 $w_t - x_t$  if  $y_t = -1$

5

CSE 446: Machine Learning

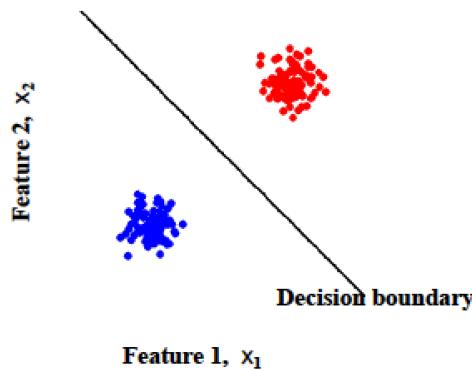
Why move like this? If  $w_t$  had BEEN  $x_t$  when  $y_t = 1$ , then we would have classified correctly.  
If  $w_t$  had BEEN  $-x_t$  when  $y_t = -1$ , then we would have classified correctly.

# How many mistakes can a perceptron make?

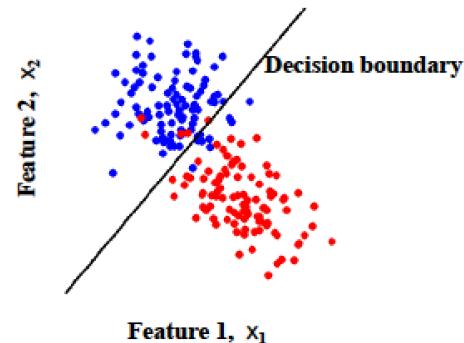
# Assumption: data is linearly separable

linearly separable: hyperplane such that the data is perfectly separated in two.

Linearly separable data



Linearly non-separable data



# Linear separability using margin

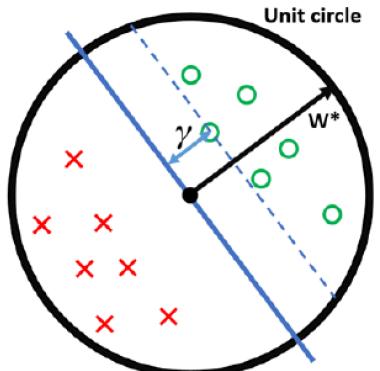


Figure by Kilian Weinberger

Data linearly separable, if there exists

- a vector  $w^*$  that defines a hyperplane (orthogonal to it) which perfectly separates the data into two groups
- a margin

such that

distance from  $x$  to hyperplane  $w^*v = 0$  is  $\geq \gamma$

THUS

if  $y_t = 1$  then  $w^* \cdot x_t \geq \gamma$

if  $y_t = -1$  then  $w^* \cdot x_t \leq -\gamma$

# Perceptron analysis: Linearly separable case

Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:  
 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Each feature vector has bounded norm:  
 $\|x\| \leq R$
- If dataset is linearly separable:

and there exists a  $\gamma > 0$  such that there exists a  $w^*$  with  $\|w^*\| = 1$  such that  $|x_i \cdot w^*| \geq \gamma$  for all  $i$

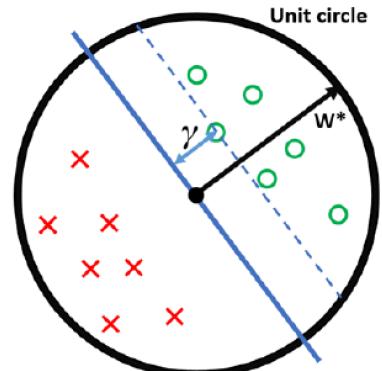


Figure by Kilian Weinberger

Then the # mistakes made by the online perceptron on this sequence is bounded by

$$\text{mistakes} \leq R^2 / \gamma^2$$

# Perceptron proof for linearly separable case

- Every time we make a mistake, we get  $w$  closer to  $w^*$ :

- Mistake at time  $t$ :
  - Taking dot product with  $w^*$ :
  - Thus after  $m$  mistakes:
- $$w^{* \{t+1\}} = w^* w_{\{t+1\}} = w^* w_t + y_t w^* x_t = w^* w_t + \gamma$$
- $w^* w_{\{t+1\}} \geq \gamma * m$

$$w_{t+1} := w_t + y_t x_t$$

$$\begin{aligned} w^* w_{\{t+1\}} &= w^* w_t + y_t w^* x_t \\ &\geq w^* w_t + \gamma \end{aligned}$$

- On the other hand, norm of  $w_{t+1}$  doesn't grow too fast:

$$\|w_{t+1}\|^2 = \|w_t\|^2 + 2y_t(w_t \cdot x_t) + \|x_t\|^2$$

- Thus, after  $m$  mistakes:

$$\|w_{\{t+1\}}\| \leq m * R^2$$

$\leq 0$  as we misclassified

- Putting all together:

$$\gamma * m \leq \|w^* w_{\{t+1\}}\| \leq \|w^*\| \|w_{\{t+1\}}\| \leq \sqrt{m} R$$

$$\gamma * m \leq \sqrt{m} R$$

$$m \leq R^2 / \gamma^2$$

$w^*$  is the perfect vector;  
defines a hyperplane  
that perfectly separates  
data into two classes.

mistake if  $y_t (x_t \cdot w^*) \leq 0$   
predicted does not match actually.

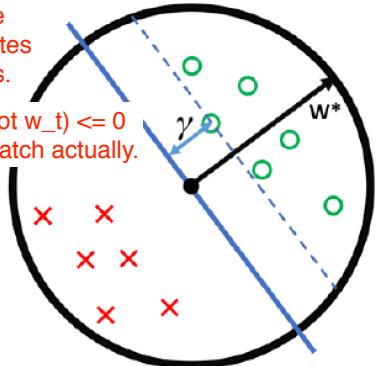
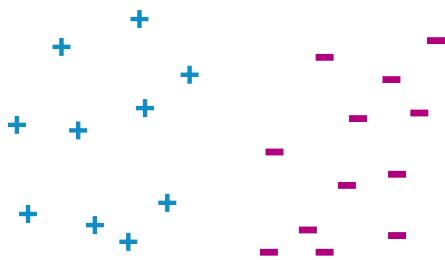


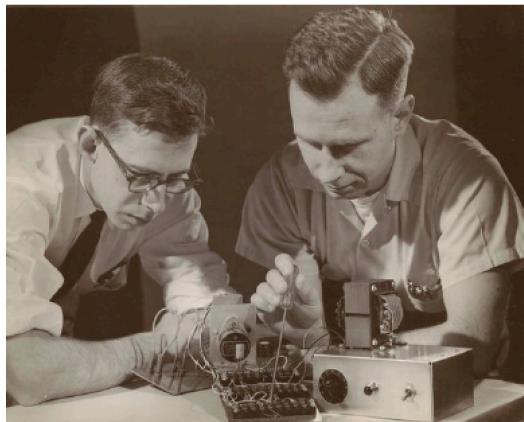
Figure by Kilian Weinberger

$\leq R^2$  by assumption that each feature has bounded norm

# Beyond linearly separable case

- Perceptron algorithm is super cool!
  - No assumption about data distribution!
    - Could be generated by an adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done forever!
    - Even if you see infinite data





Rosenblatt 1957



"the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."

*The New York Times, 1958*

## NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

WASHINGTON, July 7 (UPI)

—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

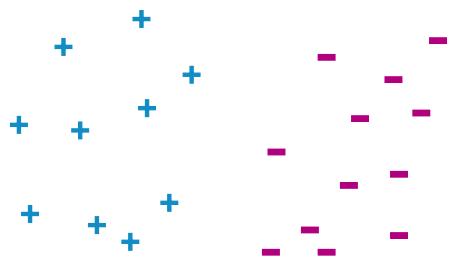
The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen..

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt. da e Learning

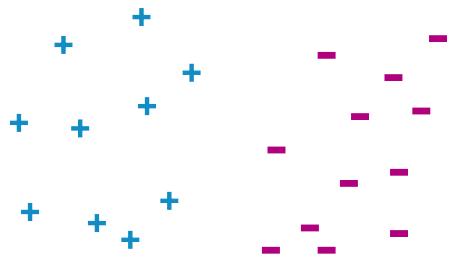
# Beyond linearly separable case

- Perceptron algorithm is super cool!
  - No assumption about data distribution!
    - Could be generated by an adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done for ever!
    - Even if you see infinite data
- Possible to motivate the perceptron algorithm as an implementation of SGD for minimizing a certain loss function.  
(see next homework)



# Beyond linearly separable case

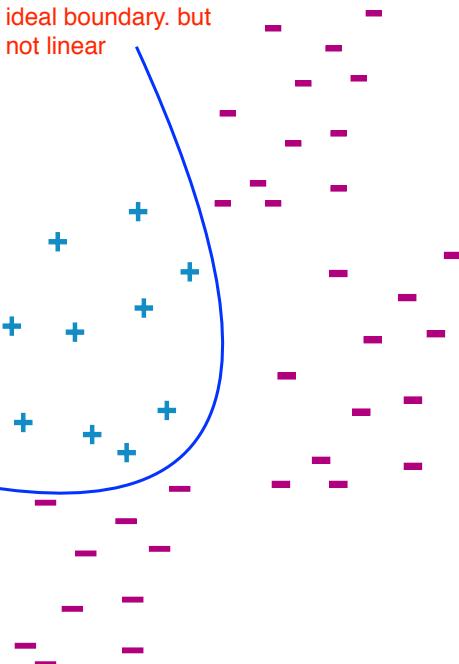
- Perceptron algorithm is super cool!
  - No assumption about data distribution!
    - Could be generated by an adversary, no need to be iid
  - Makes a fixed number of mistakes, and it's done for ever!
    - Even if you see infinite data
- However, real world is not linearly separable
  - Can't expect never to make mistakes again
  - If data not separable, cycles forever and hard to detect.
  - Even if separable, may not give good generalization accuracy (make have small margin).



i.e. small gamma... This means  
if I get future data theres a higher  
likelihood one of the points slips over the line

# The kernelized perceptron

# What to do if data are not linearly separable?



Use feature maps...

$$\phi(x) : \mathbb{R}^d \rightarrow F$$

- Start with set of inputs for each observation  $\mathbf{x} = (x[1], x[2], \dots, x[d])$  and training set:  $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$
- Define feature map that transforms each input vector  $\mathbf{x}_i$  to higher dimensional feature vector  $\phi(\mathbf{x}_i)$

Example:  $(x_i[1], x_i[2], x_i[3])^T$

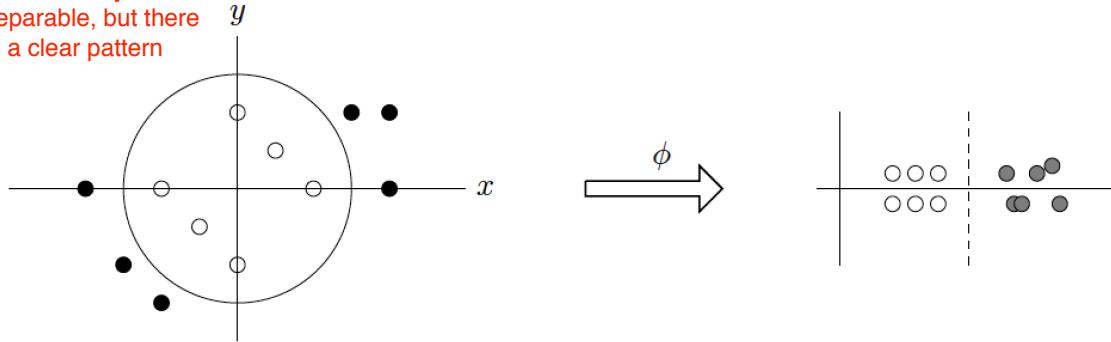
3 dimensional data. Suppose not linearly separable. Map to a higher dimensional space where it will be linearly separable.

$$\phi(\mathbf{x}_i) = (1, x_i[1], x_i[1]^2, x_i[1]x_i[2], x_i[2], x_i[2]^2, \cos(\pi x_i[3]/6))^T$$

black = 1  
white = -1

$$\vec{x} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \quad \phi(\vec{x}) = \begin{pmatrix} 1 \\ \sqrt{2} \cdot 2 \\ \sqrt{2} \cdot 3 \\ 4 \\ \sqrt{2} \cdot 6 \\ 9 \end{pmatrix}$$

NOT Linearly  
separable, but there  
is a clear pattern



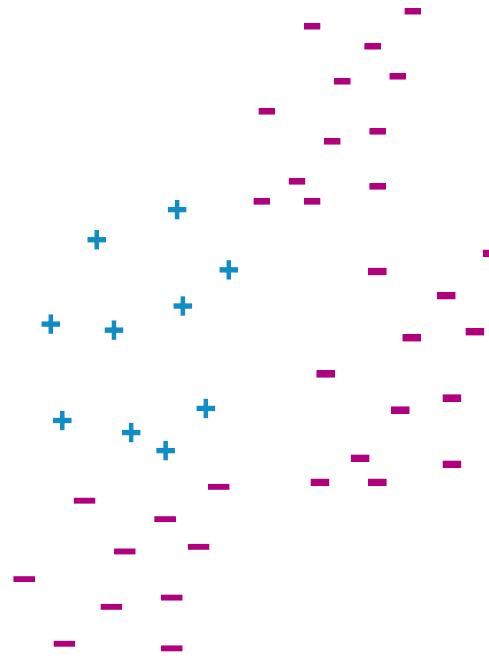
or map to  $x^2 + y^2$

Figure 5.2: Data that is not linearly separable in the input space  $\mathbb{R}^2$  but that is linearly separable in the “ $\phi$ -space,”  $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$ .

our hyperplane that linearly separates the data  
in the higher dimensional space.

For instance, the hyperplane  $\phi(\mathbf{x})^T \mathbf{w}^* = 0$  for  $\mathbf{w}^* = (-4, 0, 0, 1, 0, 1)$  circle  $x_1^2 + x_2^2 = 4$  in the original space, such as in Figure 5.2.

# What to do if data are not linearly separable?



Use feature maps...

$$\phi(x) : \mathbb{R}^d \rightarrow F$$

$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x[1] \\ x[2] \\ \vdots \\ x[1]^2 \\ x[2]^2 \\ \vdots \\ x[1]x[2] \\ x[1]x[3] \\ \vdots \end{pmatrix}$$

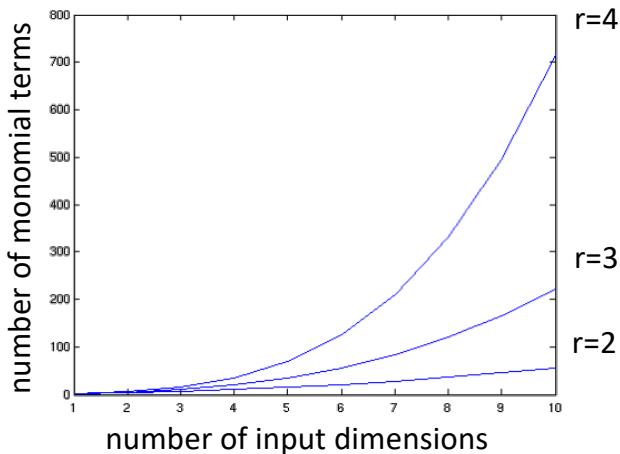
where  $\mathbf{x} = (x[1], \dots, x[d])$

$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x[1] \\ \vdots \\ \dots \\ x[1]e^{\sin(x[3])} \\ e^{-x[2]^2/\sigma^2} \\ \vdots \end{pmatrix}$$

Could even be infinite dimensional....

# Example: Higher order polynomials

Feature space can get really large really quickly!



$$\phi(\mathbf{u}) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \dots \\ u_1^2 \\ u_1 u_2 \\ \dots \\ u_1^5 u_3 u_4^4 \\ \dots \end{pmatrix}$$

d – input dimension  
r – degree of polynomial

$\binom{r+d-1}{r}$  terms

grows fast!  
 $r = 6, d = 100$   
about 1.6 billion terms

# Kernel trick

- Allows us to do these apparently intractable calculations efficiently!

# Perceptron revisited

Initialize  $\mathbf{w}_0$ .

for  $t := 1$  to  $T$

- Observe feature vector  $\mathbf{x}_t$ .
- Make a prediction:  $\hat{y} = \text{sign}(\mathbf{x}_t^T \mathbf{w}_t)$
- Observe the true label  $y_t \in \{-1, 1\}$ .
- Update the model:
  - \* If  $\hat{y} = y_t$ , then  $\mathbf{w}_{t+1} := \mathbf{w}_t$  .
  - \* Otherwise ( $\hat{y} \neq y_t$ ), then  $\mathbf{w}_{t+1} := \mathbf{w}_t + y_t \mathbf{x}_t$  .

Write weight vector in terms of mistaken data points only.

Let  $M^{(t)}$  be time steps up to  $t$  when mistakes were made:

$$\mathbf{w}_{t+1} = \sum(y_j \mathbf{x}_j) \text{ over all mistakes}$$

Prediction rule becomes:

$$\begin{aligned} & \text{sign}(\mathbf{x}_t \cdot \sum(y_j \mathbf{x}_j) \text{ for mistakes } j) \\ &= \text{sign}(\sum(y_j \mathbf{x}_t \cdot \mathbf{x}_j) \text{ for mistakes } j) \end{aligned}$$

When using high dimensional features:

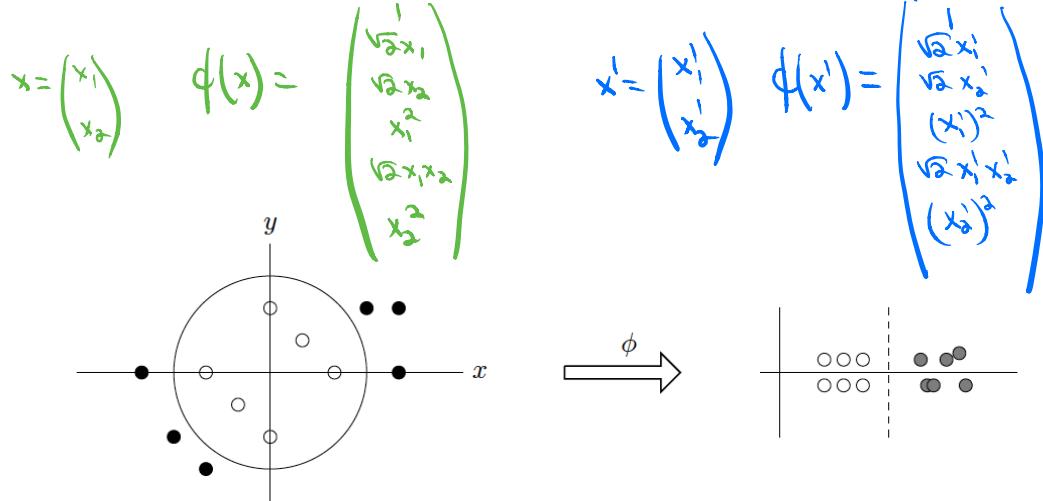
$$\text{sign}(\sum(y_j \Phi(\mathbf{x}_t) \Phi(\mathbf{x}_j)))$$

Classification only depends on inner products!

# Why does dependence on inner products help?

Because sometimes they can be computed **much much much more efficiently**.

$$\phi(\mathbf{u}) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \dots \\ u_1^2 \\ u_1 u_2 \\ \dots \\ u_1^5 u_3 u_{17}^4 \\ \dots \end{pmatrix} \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \dots \\ u_1^2 \\ u_1 u_2 \\ \dots \\ u_1^5 u_3 u_{17}^4 \\ \dots \end{pmatrix} \cdot \begin{pmatrix} 1 \\ v_1 \\ v_2 \\ \dots \\ v_1^2 \\ v_1 v_2 \\ \dots \\ v_1^5 v_3 v_{17}^4 \\ \dots \end{pmatrix}$$
$$\binom{r+d-1}{r}$$



**Figure 5.2:** Data that is not linearly separable in the input space  $\mathbb{R}^2$  but that is linearly separable in the “ $\phi$ -space,”  $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$ , corresponding to the kernel function  $K(\mathbf{x}, \mathbf{x}') = (1 + x_1x'_1 + x_2x'_2)^2$ .

For instance, the hyperplane  $\phi(\mathbf{x})^T \mathbf{w}^* = 0$  for  $\mathbf{w}^* = (-4, 0, 0, 1, 0, 1)$  circle  $x_1^2 + x_2^2 = 4$  in the original space, such as in Figure 5.2.

# Dot-product of polynomials: why useful?

- Because sometimes they can be computed much much much more efficiently.

$$\phi(\mathbf{u}) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \vdots \\ u_1^2 \\ u_1 u_2 \\ \vdots \\ u_1^5 u_3 u_{17}^4 \\ \vdots \end{pmatrix}$$
$$\phi(u) \cdot \phi(v) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \vdots \\ u_1^2 \\ u_1 u_2 \\ \vdots \\ u_1^5 u_3 u_{17}^4 \\ \vdots \end{pmatrix} \cdot \begin{pmatrix} 1 \\ v_1 \\ v_2 \\ \vdots \\ v_1^2 \\ v_1 v_2 \\ \vdots \\ v_1^5 v_3 v_{17}^4 \\ \vdots \end{pmatrix}$$

So with appropriate constants in front

u and v are only dimension d! MUCH smaller

$$\phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = (1 + \mathbf{u} \cdot \mathbf{v})^r$$

$$(1 + \mathbf{u} \cdot \mathbf{v})^r = [1 + u_1 v_1 + u_2 v_2 + \dots + u_d v_d]^r$$

Typical term:  $C \cdot (u_1 v_1)^{i_1} (u_2 v_2)^{i_2} \cdots (u_d v_d)^{i_d}$  where  $\sum_{j=1}^d i_j \leq r$

constants are unimportant...

# Kernel trick

- Allows us to do these apparently intractable calculations efficiently!
- If we can compute inner products efficiently, can run the algorithm (e.g., Perceptron) efficiently.
- Hugely important idea in ML.

# Finally, the kernel trick

## Kernelized perceptron

A Kernel function  $K(x, x')$  is a function that implements inner products in the feature space

$K(x, x') = \phi(x) \cdot \phi(x')$   
where  $\phi: \mathbb{R}^d \rightarrow \mathcal{F}$

where  $\phi$  maps from the features to some higher dimensional feature space.

- Every time you make a mistake, remember  $(x_t, y_t)$
- Kernelized perceptron prediction for  $x$ :

$$\begin{aligned}\text{sign}(\mathbf{w}_t \cdot \phi(\mathbf{x})) &= \sum_{i \in M^{(t)}} y_i (\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x})) \\ &= \sum_{i \in M^{(t)}} y_i K(\mathbf{x}_i, \mathbf{x}).\end{aligned}$$

$$K(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = (1 + \mathbf{u} \cdot \mathbf{v})^r$$

# Common kernels

- Polynomials of degree exactly p

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^p$$

- Polynomials of degree up to p

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^p$$

- **Gaussian (squared exponential) kernel**

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

- Many many many others

# Kernelizing a machine learning algorithm

- Prove that the solution is always in span of training points. I.e.,
- Rewrite the algorithm so that all training or test inputs are accessed only through inner products with other inputs.
- Choose (or define) a kernel function and substitute  $K(\mathbf{x}_i, \mathbf{x}_j)$  for  $\mathbf{x}_i^T \mathbf{x}_j$

So with appropriate constants in front in definition of  $\phi(\cdot)$

$$K(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u}) \cdot \phi(\mathbf{v}) = (1 + \mathbf{u} \cdot \mathbf{v})^r$$

# What you need to know

- Linear separability in higher-dim feature space
- The kernel trick
- Kernelized perceptron
- Polynomial and other common kernels (will see more later)

# Support Vector Machines

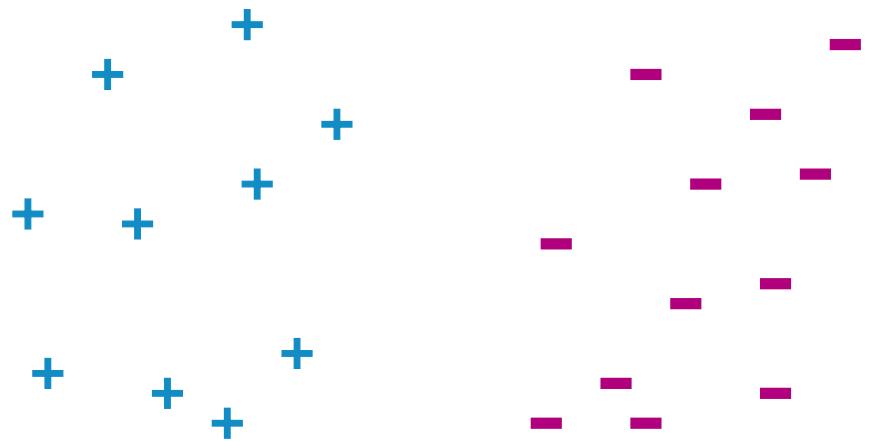
CSE 446: Machine Learning

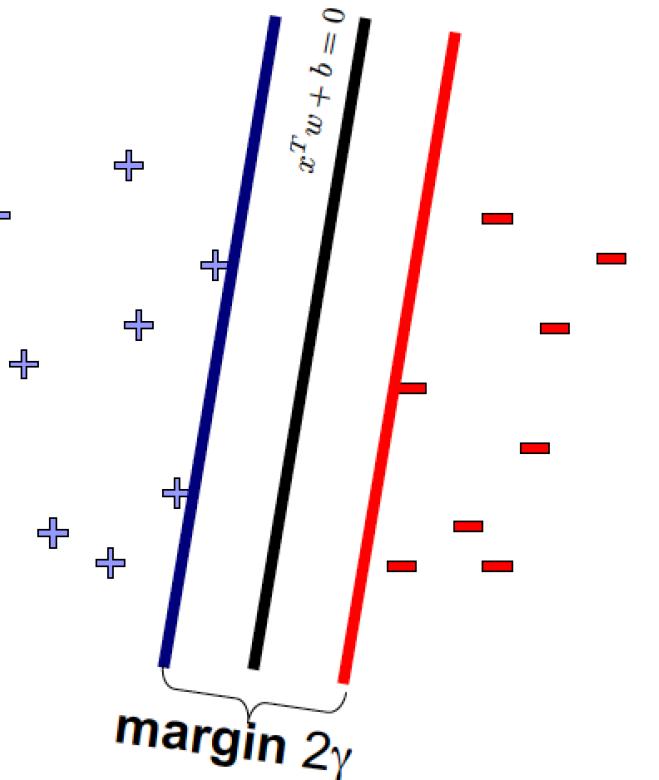
Slides by Emily Fox + Kevin Jamieson + others

Presented by Anna Karlin

May 1, 2019

# Linear classifiers—Which line is better?





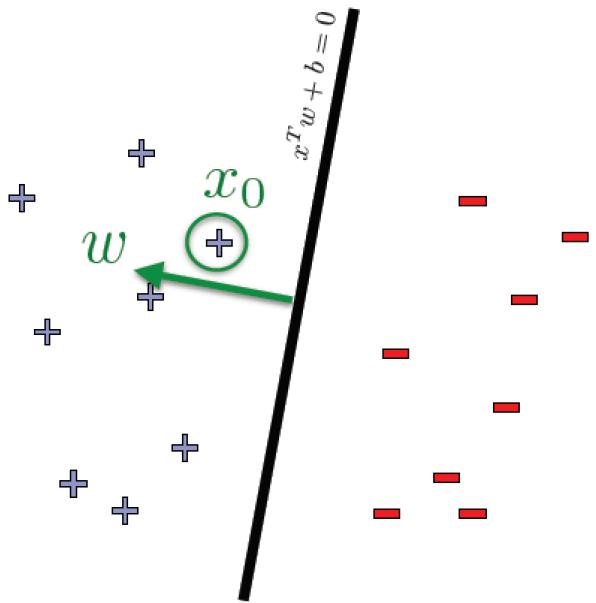
©2018 Kevin Jamieson

©2017 Emily Fox

CSE 446: Machine Learning

# Maximizing the margin for linearly separable data

# Given hyperplane, what is margin?

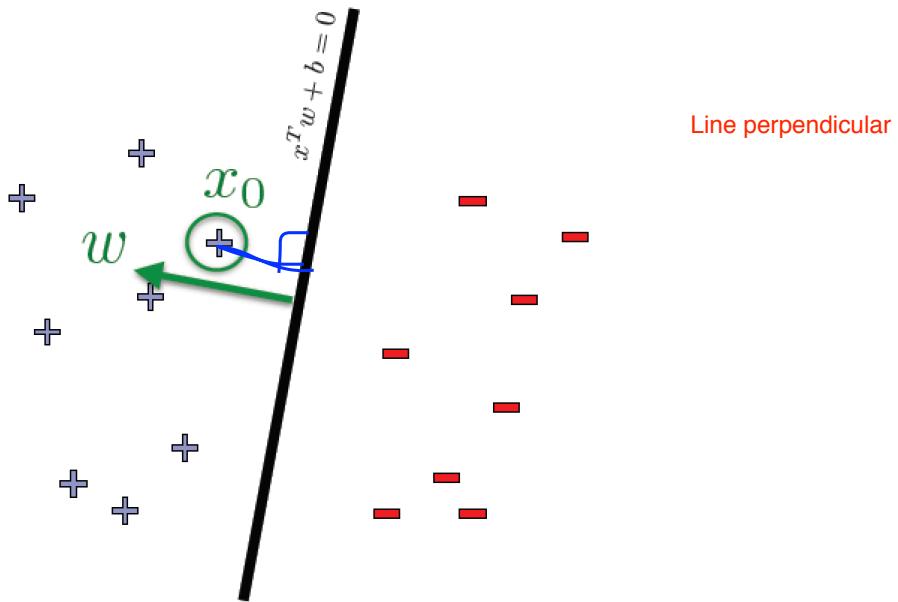


Distance from  $x_0$  to  
hyperplane defined  
by  $x^T w + b = 0$ ?

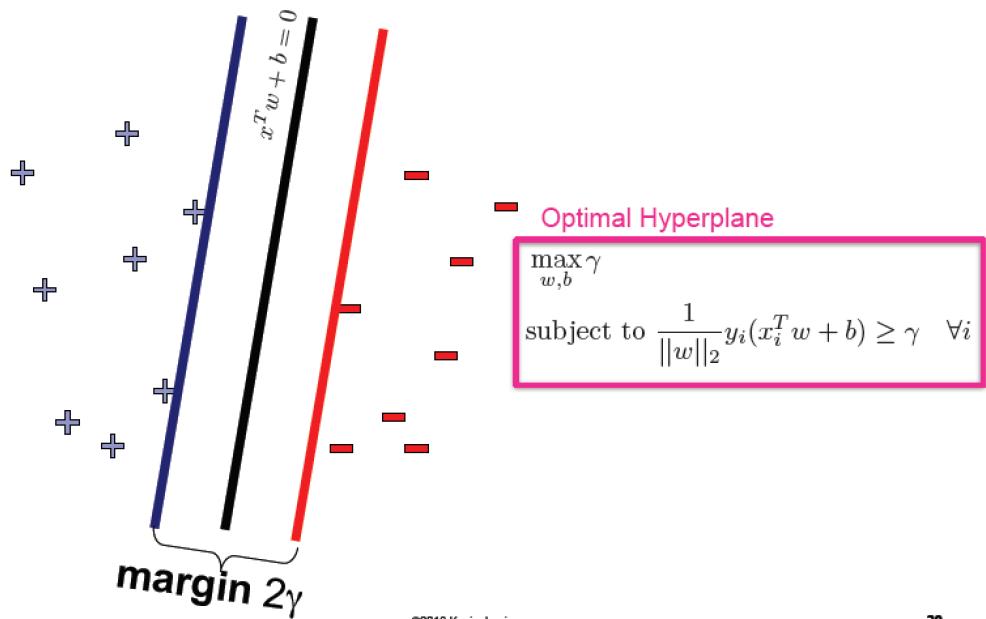
larger margin between hyperplane and data = larger margin of error  
in forming a classification.

Choose hyperplane that is as far away from the nearest data point

# Given hyperplane, what is margin?



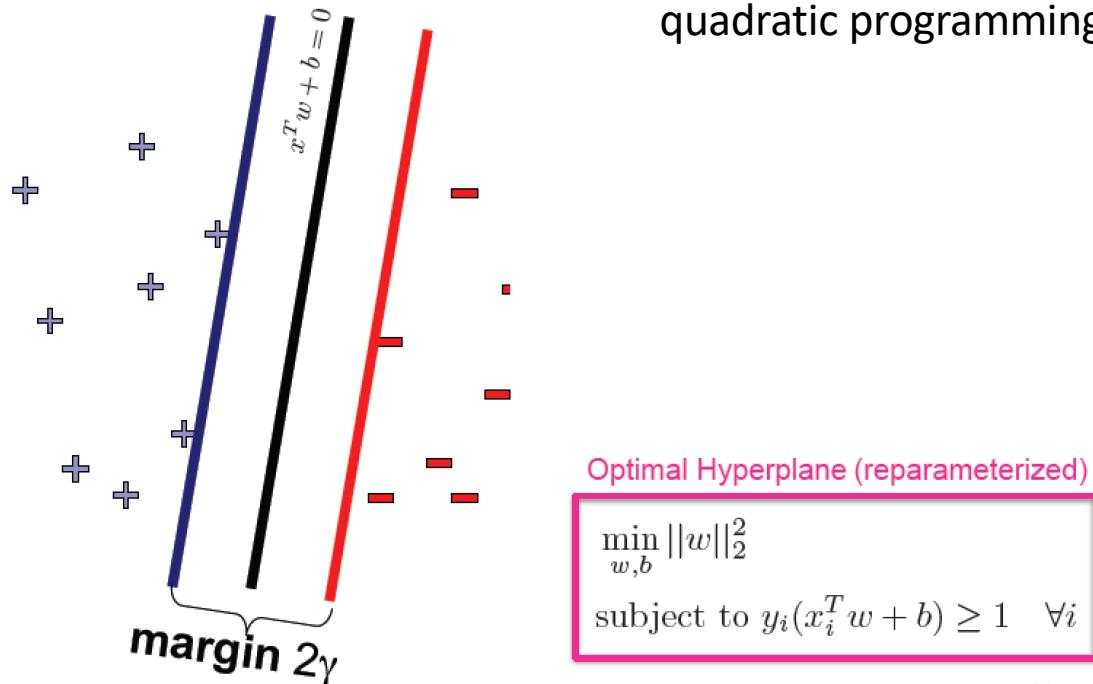
# Our optimization problem



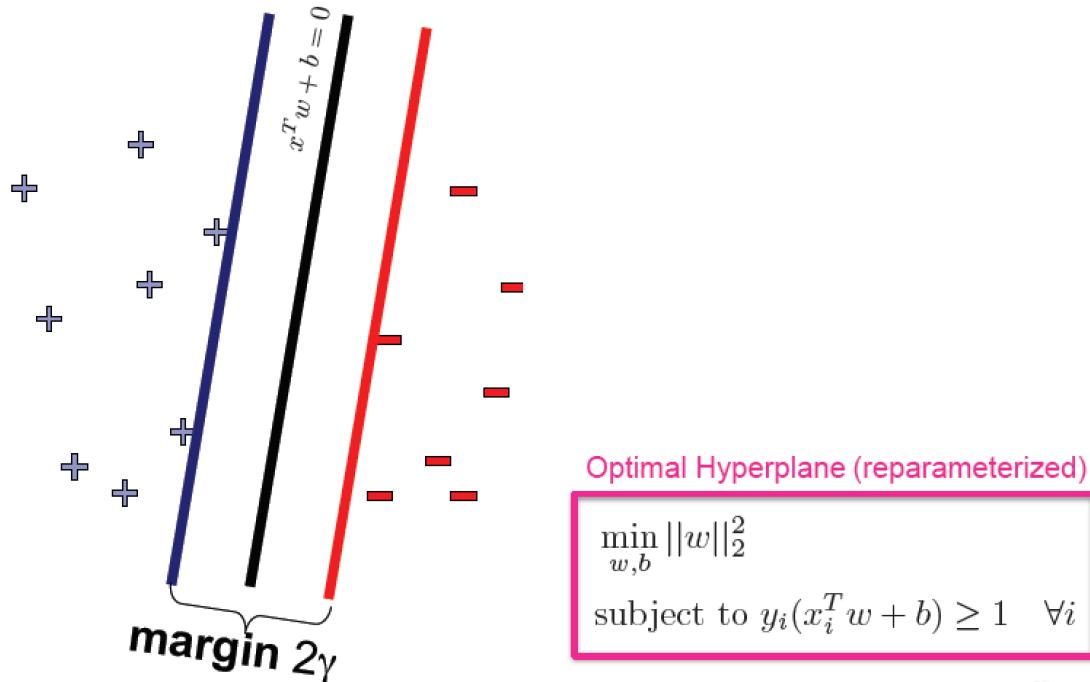
# END LECTURE WED MAY 1

## Final version

Solvable efficiently –  
quadratic programming problem



# What are support vectors?



# What if the data are not linearly separable?

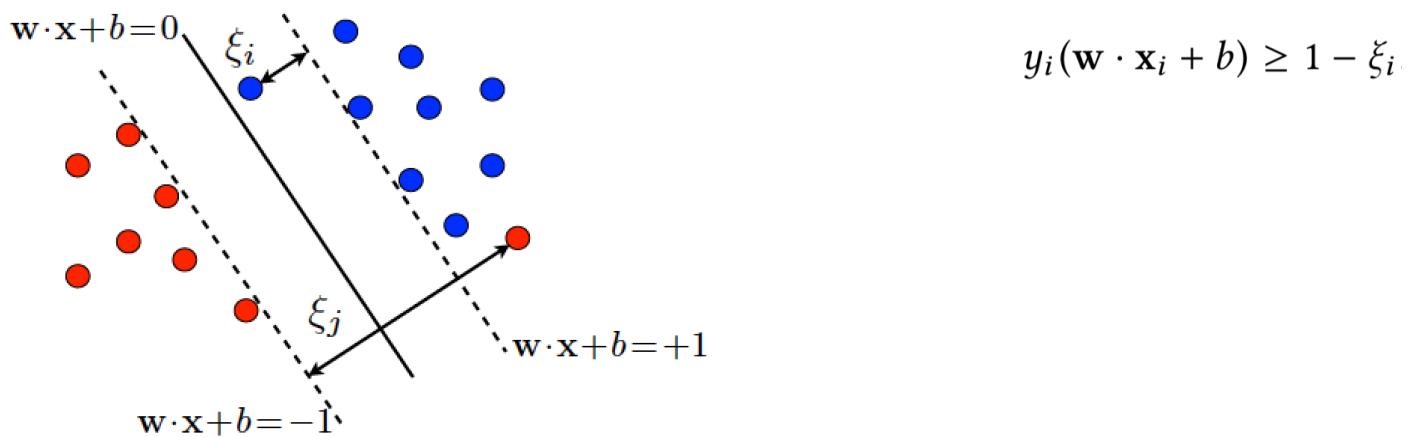
# What if data are not linearly separable?

Use feature maps...



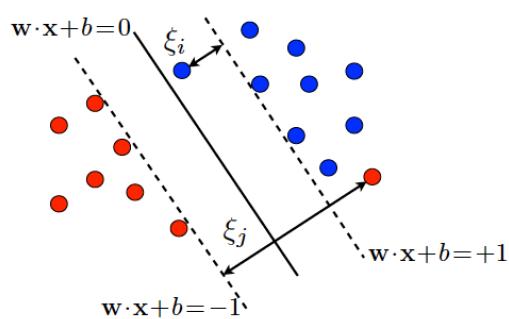
SVMs can be kernelized!!!!

# What if data are still not linearly separable?



Courtesy Mehryar Mohri

# What if data are still not linearly separable?



$$\min_{\mathbf{w}, b, \xi} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i$$
$$\xi_i \geq 0 \quad \forall i.$$

Courtesy Mehryar Mohri

# Final objective

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i((\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2$$

# Gradient descent for SVMs

# Minimizing regularized hinge loss (aka SVMs)

- Given a dataset:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
- Minimize regularized hinge loss:  $\frac{1}{n} \sum_i (1 - y_i((\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2)$

# Subgradient of hinge loss

- Hinge loss:  $\ell((\mathbf{x}, y), \mathbf{w}) = (1 - y(\mathbf{w}^T \mathbf{x} + b))_+$

- Subgradient of hinge loss:

$$\partial_{\mathbf{w}} \ell((\mathbf{x}, y), \mathbf{w}) = \begin{cases} \cdot & y(\mathbf{w}^T \mathbf{x} + b) > 1 \\ \cdot & y(\mathbf{w}^T \mathbf{x} + b) < 1 \\ \cdot & y(\mathbf{w}^T \mathbf{x} + b) = 1 \end{cases}$$

# Subgradient of hinge loss

- Hinge loss:  $\ell((\mathbf{x}, y), \mathbf{w}) = (1 - y(\mathbf{w}^T \mathbf{x} + b))_+$

- Subgradient of hinge loss:

$$\partial_{\mathbf{w}} \ell((\mathbf{x}, y), \mathbf{w}) = \begin{cases} \mathbf{0} & y(\mathbf{w}^T \mathbf{x} + b) > 1 \\ -y\mathbf{x} & y(\mathbf{w}^T \mathbf{x} + b) < 1 \\ [-y\mathbf{x}, \mathbf{0}] & y(\mathbf{w}^T \mathbf{x} + b) = 1 \end{cases}$$

In one line:

$$\partial_{\mathbf{w}} \ell((\mathbf{x}, y), \mathbf{w}) = \mathbb{I}\{y(\mathbf{w}^T \mathbf{x} + b) \leq 1\}(-y\mathbf{x})$$

# Subgradient descent for hinge minimization

- Given data:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$

- Want to minimize:

$$\frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{n} \sum_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2$$

- As we've discussed, subgradient descent works like gradient descent:
  - But if there are multiple subgradients at a point, just pick (any) one:

$$\begin{aligned}\mathbf{w}_{t+1} &:= \mathbf{w}_t - \eta \left( \frac{1}{n} \sum_{i=1}^n \partial_{\mathbf{w}} \ell((\mathbf{x}_i, y_i), \mathbf{w}) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t - \eta \left( \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\} (-y_i \mathbf{x}_i) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t + \eta \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\} (y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t.\end{aligned}$$

# SVM

- **Gradient Descent Update**

$$\begin{aligned}\mathbf{w}_{t+1} &:= \mathbf{w}_t - \eta \left( \frac{1}{n} \sum_{i=1}^n \partial_{\mathbf{w}} \ell((\mathbf{x}_i, y_i), \mathbf{w}) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t - \eta \left( \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\}(-y_i \mathbf{x}_i) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t + \eta \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\}(y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t.\end{aligned}$$

- **SGD update**

$$\mathbf{w}_{t+1} := \mathbf{w}_t + \eta \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\}(y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t.$$

# Machine learning problems

- Given i.i.d. data set:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$$

- Find parameters  $\mathbf{w}$  to minimize average loss  
(or regularized version):

$$\frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

Squared loss:

$$\ell_i(\mathbf{w}) = (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

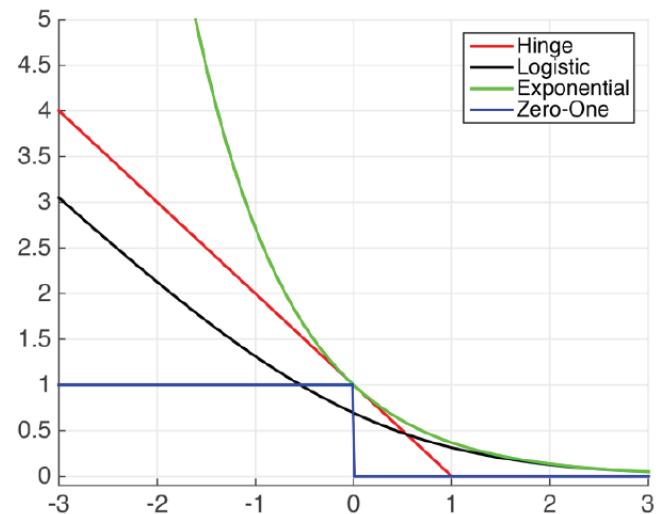
Logistic loss:

$$\ell_i(\mathbf{w}) = \log(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w}))$$

Hinge loss:

$$\ell_i(\mathbf{w}) = \max\{0, 1 - y_i \mathbf{x}_i^T \mathbf{w}\}$$

Courtesy Killian Weinberger



# What you need to know...

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
  - Hinge loss
  - a.k.a. adding slack variables
- Can optimize SVMs with SGD
  - Many other approaches possible