

- additional /shuffled office hours for upcoming week
see web site
- $\frac{1}{2}$ g you will take midterm in a different room

The Perceptron Algorithm

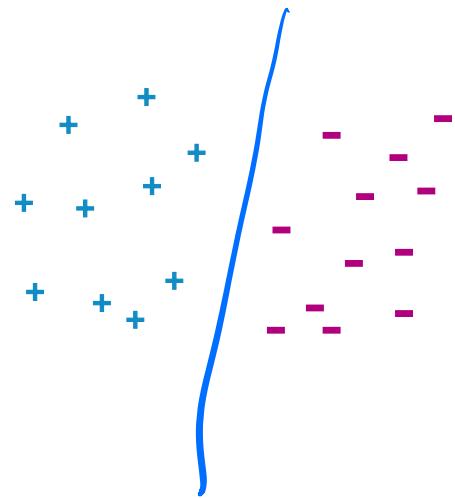
CSE 446: Machine Learning
Slides by Emily Fox (mostly)
Presented by Anna Karlin

April 29, 2019

Can we do binary classification without a model?

$$\Pr(Y=y | x, \omega) = \frac{1}{1 + \exp(-y \omega^\top x)}$$

$$\{x_i, y_i\} \quad x_i \in \mathbb{R}^d \quad y_i \in \{-1, 1\}$$



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 Find \vec{w}, b

- Prediction:

given \vec{x} $y = \begin{cases} +1 & \vec{w}^\top \vec{x} + b > 0 \\ -1 & \text{o.w.} \end{cases}$

- Training:

- Initialize weight vector: $\vec{w}_0 := \vec{0}$

- At each time step:

- Observe features: \vec{x}_+
- Make prediction: $\hat{y}_+ = \text{sign}(\vec{w}_+^\top \vec{x}_+ + b)$
- Observe true class: y_+

- Update model:

- If prediction is not equal to truth

$$\vec{x} \in \mathbb{R}^d$$
$$\vec{w} \in \mathbb{R}^d$$

$$\text{sign}(\vec{w}^\top \vec{x} + b)$$



$$y \neq \hat{y}_+$$

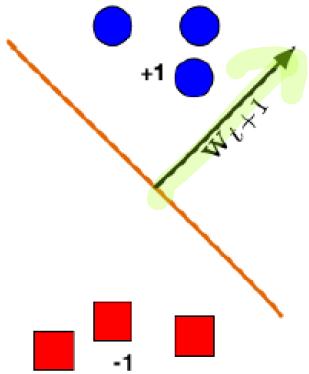
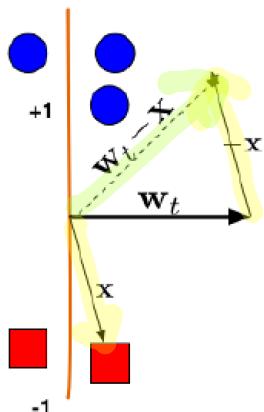
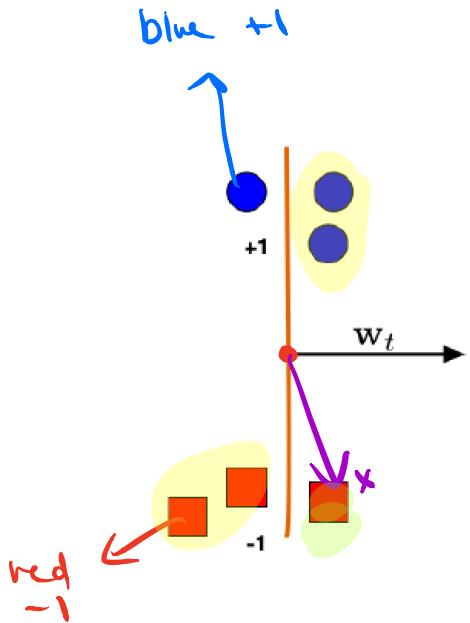
else

$$\vec{w}_{+1} = \vec{w}_+ + y_+ \vec{x}_+$$

$$\vec{w}_{+1} := \vec{w}_+$$

online setting
or batch.

Picture due to Killian Weinberger



$$x_t \in \mathbb{R}^d$$

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

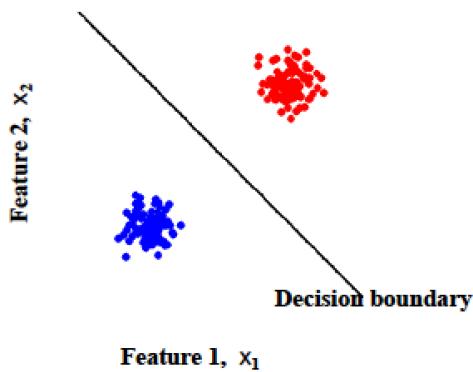
$$w_{t+1} = \begin{cases} w_t + x_t & \\ w_t - x_t & \end{cases}$$

$$\begin{aligned} y_t &= 1 \\ y_t &= -1 \end{aligned}$$

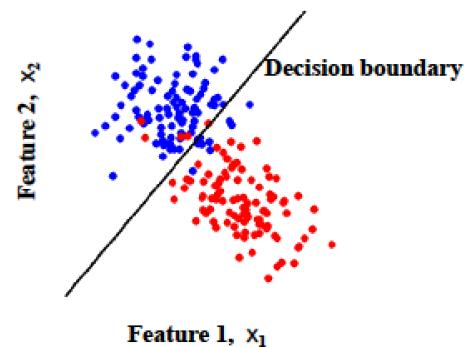
How many mistakes can a perceptron make?

Assumption: data is linearly separable

Linearly separable data



Linearly non-separable data



Linear separability using margin

$$x_i \cdot w^* = \|x_i\| \|w^*\| \cos \theta$$

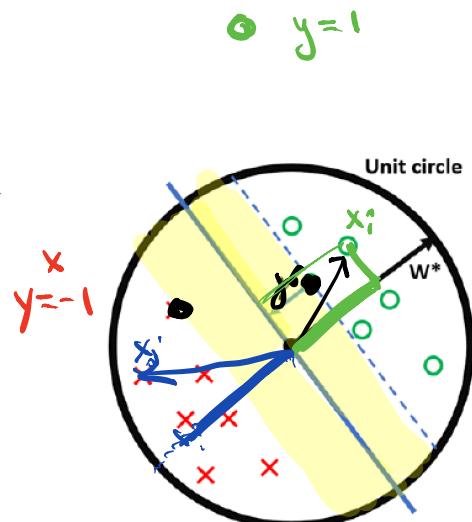


Figure by Kilian Weinberger

Data linearly separable, if there exists

- a vector w^* $\|w^*\|=1$
- a margin γ s.t. dist from pt x_i to hyperplane ($w^* \cdot x = 0$)

such that

$$\begin{cases} y_i = 1 & x_i \cdot w^* \geq \gamma \\ y_i = -1 & x_i \cdot w^* \leq -\gamma \end{cases}$$

©2017 Emily Fox

$$y_i (x_i \cdot w^*) \geq \gamma$$

6

Perceptron analysis: Linearly separable case

Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:
- Each feature vector has bounded norm:

$$\|x_i\| \leq R \quad \text{and } \gamma > 0 \text{ s.t. } y_i(x_i; w^*) \geq \gamma \quad \forall i$$

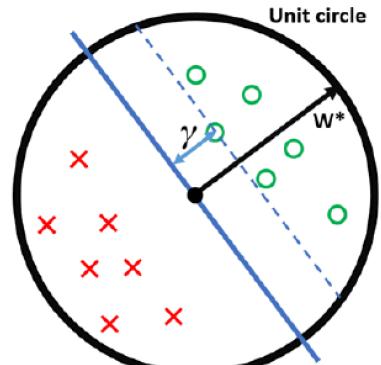


Figure by Kilian Weinberger

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

$$\frac{R^2}{\gamma \alpha}$$

$$\|x_i\| \leq R \Leftrightarrow$$

and $\gamma > 0$ s.t. $y_i(x_i \cdot w^*) \geq \gamma \quad \forall i$

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:

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

$$y_t(x_t \cdot w_t) \leq 0$$

$$\begin{aligned} w_{t+1} \cdot w^* &= w^* \cdot (w_t + y_t x_t) \\ &= w^* \cdot w_t + y_t w^* \cdot x_t \\ &\geq 0 \end{aligned}$$

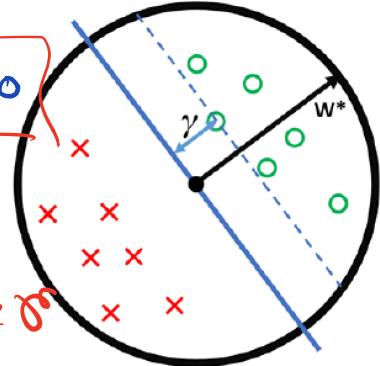


Figure by Kilian Weinberger

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

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t\|^2 + 2y_t(w_t \cdot x_t) + \|x_t\|^2 \\ (w_t + y_t x_t) \cdot (w_t + y_t x_t) &= \|w_t\|^2 + 2y_t(w_t \cdot x_t) + \|x_t\|^2 \\ - \text{Thus, after } m \text{ mistakes:} & \leq 0 \leq R^2 \end{aligned}$$

$$\|w_{t+1}\|^2 \leq mR^2$$

- Putting all together:

$$\begin{aligned} \gamma m &\leq w^* \cdot w_{t+1} \leq \|w^*\| \|w_{t+1}\| \leq \sqrt{m} R \end{aligned}$$

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

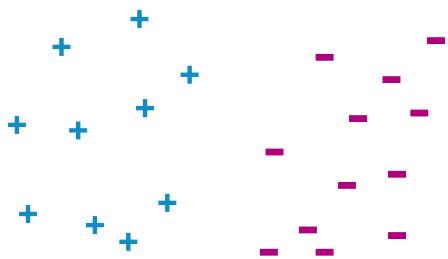
$$\sqrt{m} \leq R$$

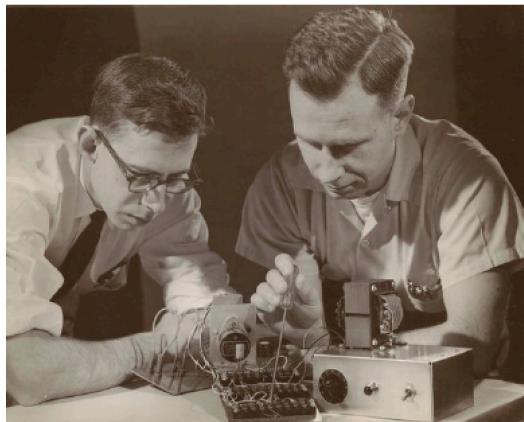
$$m \leq \frac{R^2}{\gamma^2}$$



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

X O

δ X

©2017 Emily Fox

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.

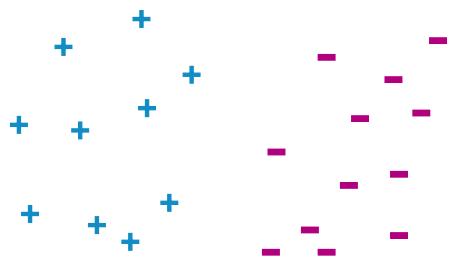
The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen..

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

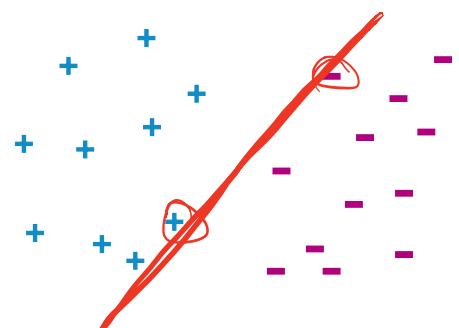
Beyond linearly separable case

- 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)



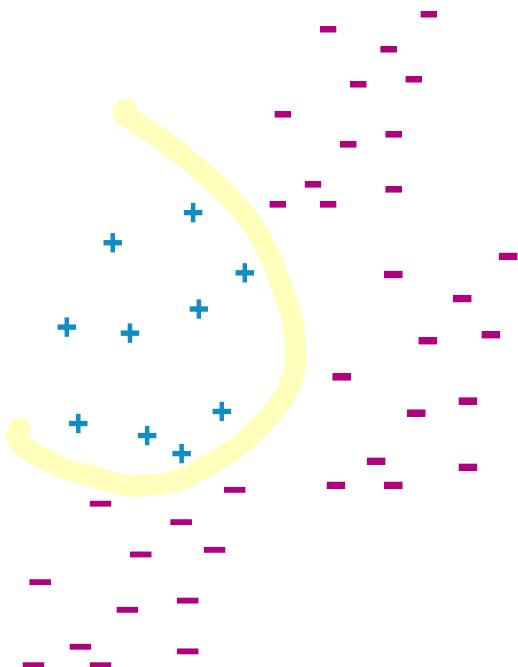
Beyond linearly separable case

- 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).



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: $(x_i[1], x_i[2], x_i[3])^\top$

$$h(x)$$

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

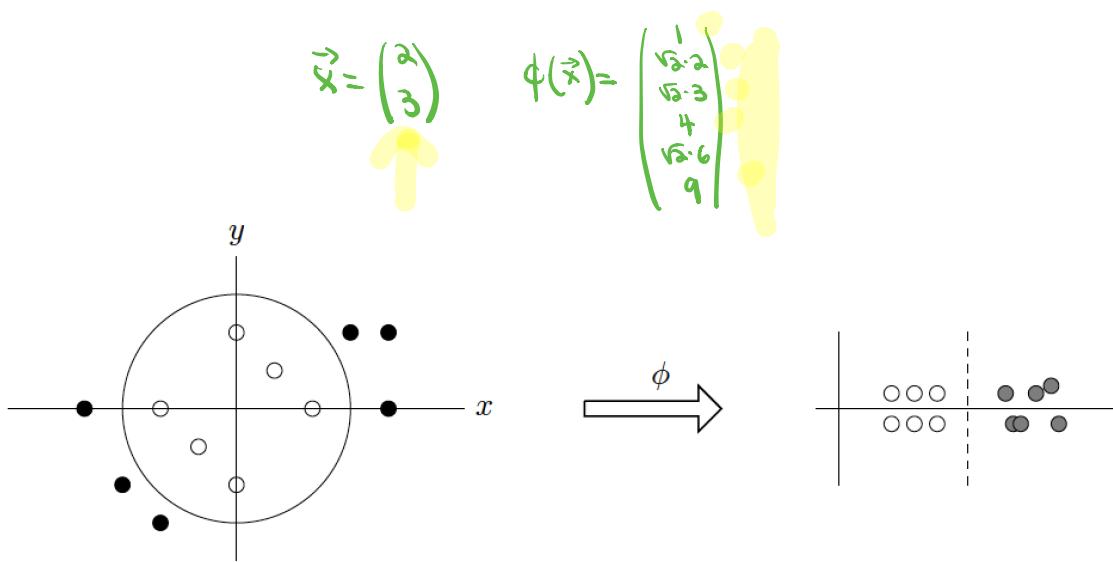


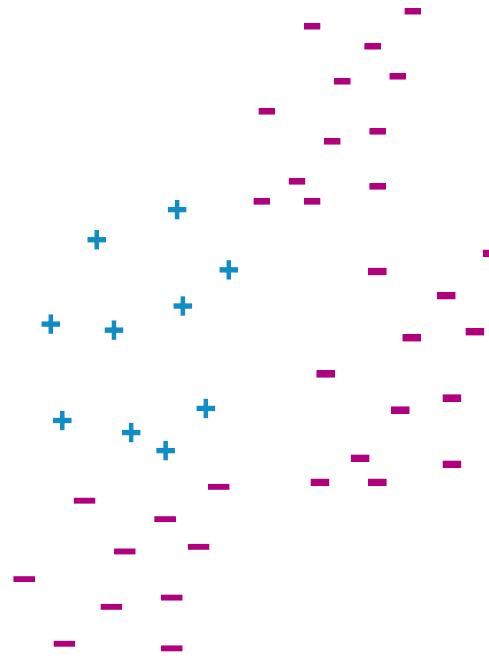
Figure 5.2: Data that is not linearly separable in the input space \mathbb{R}^2 but that is linearly separable in the “ ϕ -space,” $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$.

$$\begin{aligned} & \text{---} \\ & -4 + x_1^2 + x_2^2 = 0 \end{aligned}$$

$$x_1^2 + x_2^2 = 4$$

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...

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

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

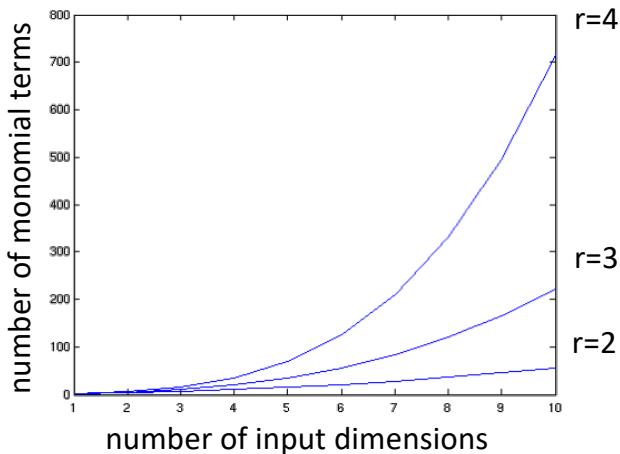
where $\mathbf{x} = (x[1], \dots, x[d])$

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

Could even be infinite dimensional....

Example: Higher order polynomials

Feature space can get really large really quickly!



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

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

$$\mathbf{x} \rightarrow \phi(\mathbf{x})$$

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 $\mathbf{w}_{t+1} := \mathbf{w}_t + y_t \mathbf{x}_t$.

When using high dimensional features:

$$\text{sign}\left(\sum_{j \in M^{(t)}} y_j \underbrace{\phi(x_j) \cdot \phi(x_j)}_{\text{inner product}}\right)$$

Classification only depends on inner products!

Write weight vector in terms of mistaken data points only.

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

$$\mathbf{w}_{t+1} = \sum_{j \in M^{(t)}} y_j \mathbf{x}_j$$

Prediction rule becomes:

$$\begin{aligned} & \text{sign}\left(\mathbf{x}_+ \cdot \left(\sum_{j \in M^{(t)}} y_j \mathbf{x}_j\right)\right) \\ &= \text{sign}\left(\sum_{j \in M^{(t)}} y_j \mathbf{x}_+ \cdot \mathbf{x}_j\right) \end{aligned}$$

Why does dependence on inner products help?

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

$$\phi(\mathbf{u}) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \dots \\ u_1^2 \\ u_1 u_2 \\ \dots \\ u_1^5 u_3 u_{17}^4 \\ \dots \end{pmatrix} \quad \phi(u) \cdot \phi(v) = \underbrace{\dots}_{\text{red underline}} \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ \dots \\ u_1^2 \\ u_1 u_2 \\ \dots \\ u_1^5 u_3 u_{17}^4 \\ \dots \end{pmatrix} \cdot \begin{pmatrix} 1 \\ v_1 \\ v_2 \\ \dots \\ v_1^2 \\ v_1 v_2 \\ \dots \\ v_1^5 v_3 v_{17}^4 \\ \dots \end{pmatrix}$$

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

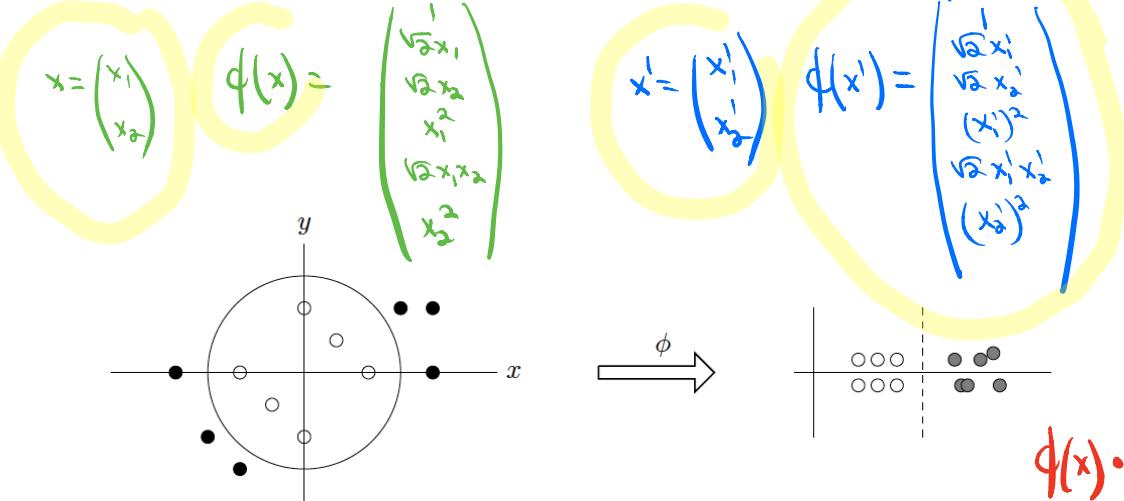


Figure 5.2: Data that is not linearly separable in the input space \mathbb{R}^2 but that is linearly separable in the “ ϕ -space,” $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$, corresponding to the kernel function $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}_1 \cdot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

$$(1 + \mathbf{x}_1 \cdot \mathbf{x}'_1 + \mathbf{x}_2 \cdot \mathbf{x}'_2)^2$$

$$\begin{aligned}\phi(\mathbf{x}) \cdot \phi(\mathbf{x}') &= 1 + 2\mathbf{x}_1 \cdot \mathbf{x}'_1 \\ &\quad + 2\mathbf{x}_2 \cdot \mathbf{x}'_2 \\ &\quad + \mathbf{x}_1^2 (\mathbf{x}'_1)^2 + \dots\end{aligned}$$

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.

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

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

So with appropriate constants in front

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

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

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

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

- Every time you make a mistake, remember (\mathbf{x}_t, y_t)

It kernel fn $K(\mathbf{x}, \mathbf{x}')$
is a fn that implements
inner products in the feature space
 $\underline{K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}')}}$

$$\phi: \mathbb{R}^d \rightarrow \mathcal{F}$$

- Kernelized perceptron prediction for \mathbf{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.,

$$\omega = \sum \alpha_i x_i$$

- 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 $K(x_i, x_j)$ for $\underline{\underline{x_i^T x_j}}$

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

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

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

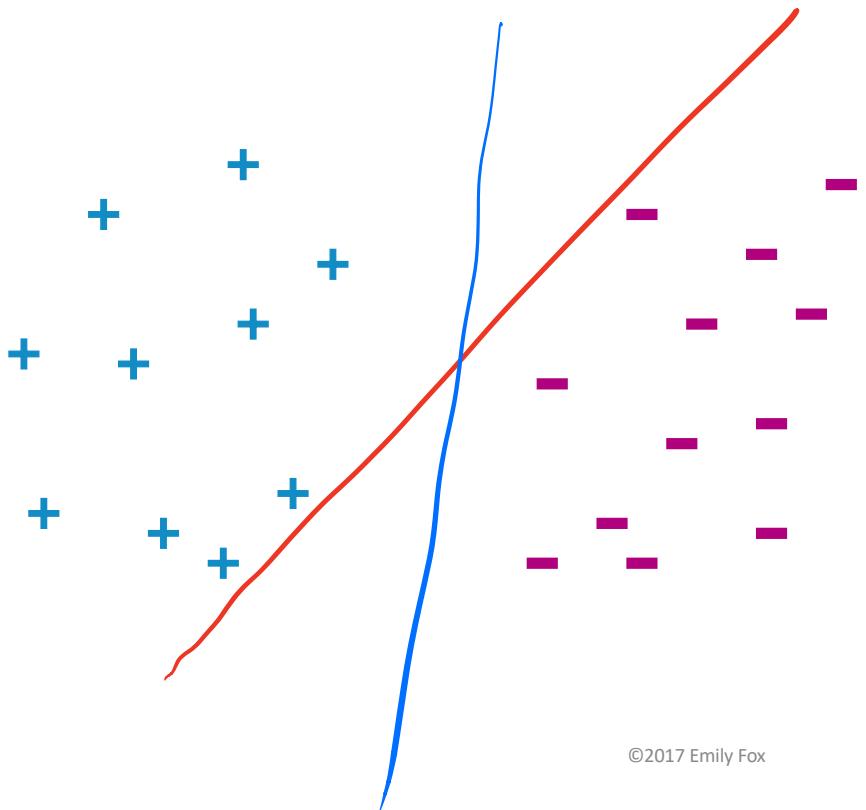
CSE 446: Machine Learning

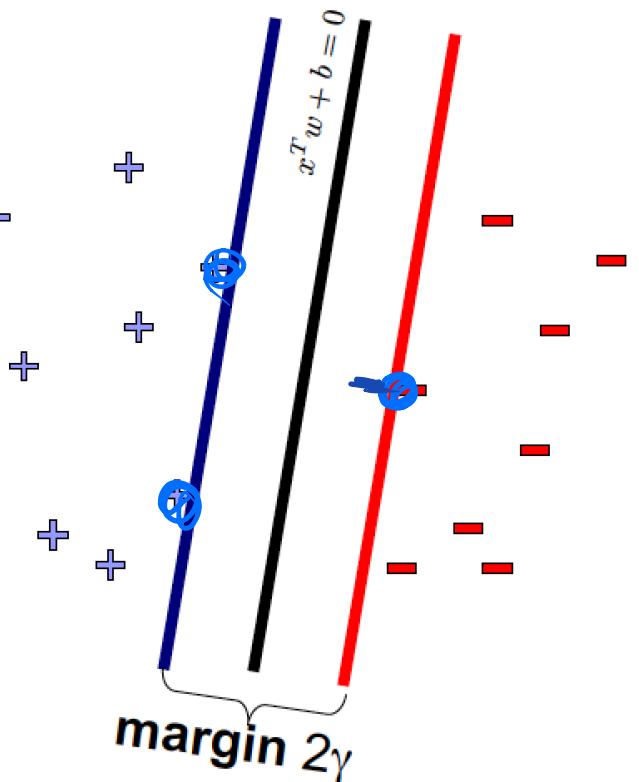
Slides by Emily Fox + Kevin Jamieson + others

Presented by Anna Karlin

May 1, 2019

Linear classifiers—Which line is better?



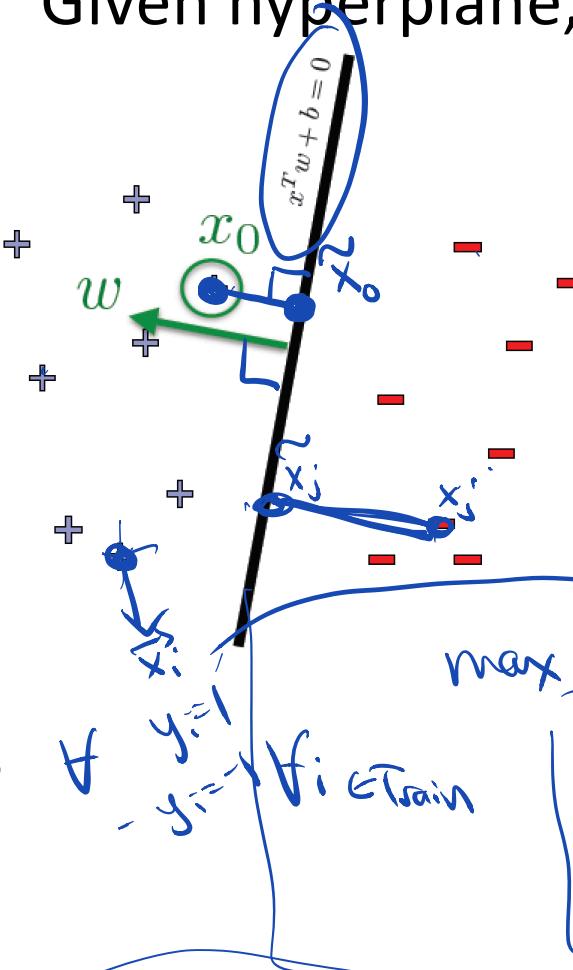


Maximize
the min dist
from the hyperplane
to an data pt β

Maximizing the margin for linearly separable data

$$H = \{x \mid \langle x, w \rangle + b = 0\}$$

Given hyperplane, what is margin?



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

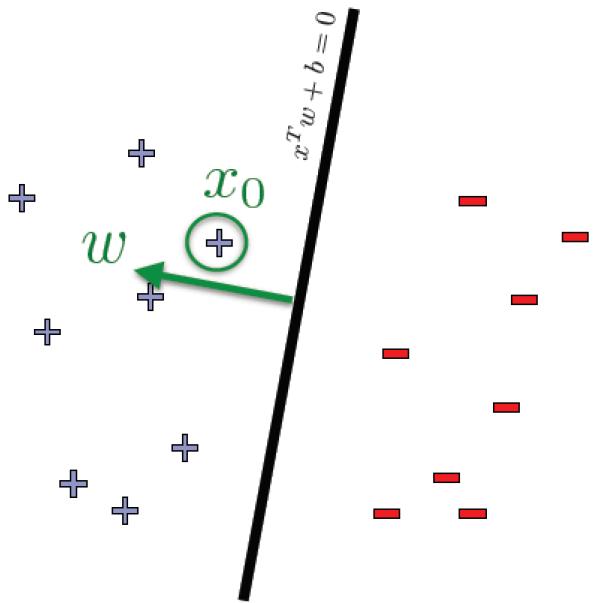
$$\begin{aligned} \|x_0 - \tilde{x}_0\| &= \frac{(x_0 - \tilde{x}_0) \cdot w}{\|w\|} \\ &= \frac{w \cdot x_0 + b}{\|w\|} \end{aligned}$$

max γ^*

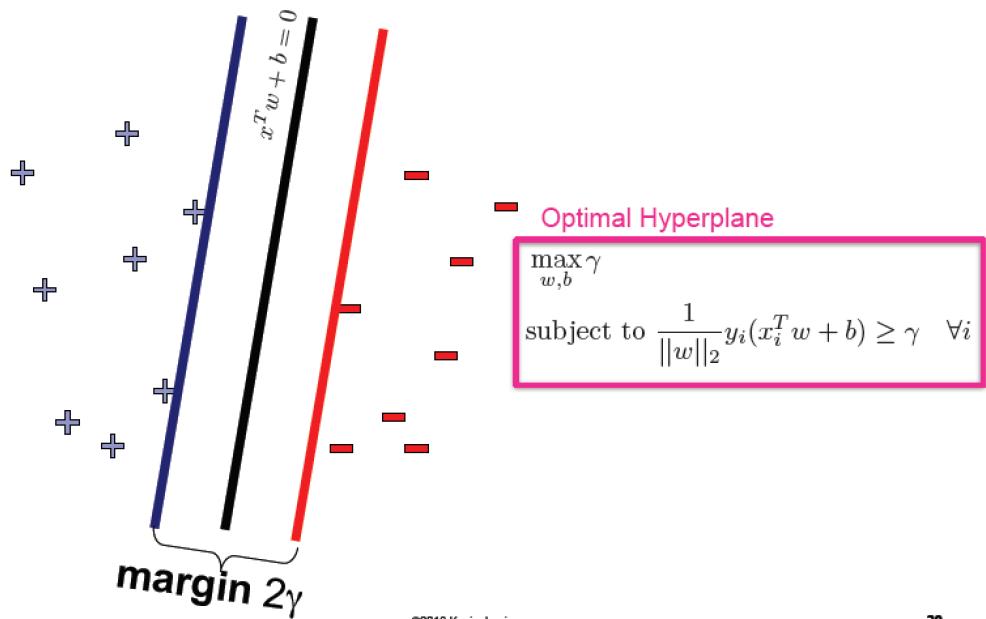
$$\left| \frac{y_i(w \cdot x_i + b)}{\|w\|} \right| \geq \gamma^*$$

$$(x_j - \tilde{x}_j) \cdot \frac{w}{\|w\|} = \|x_j - \tilde{x}_j\| \cos \theta = -\|x_j - \tilde{x}_j\|$$

Given hyperplane, what is margin?

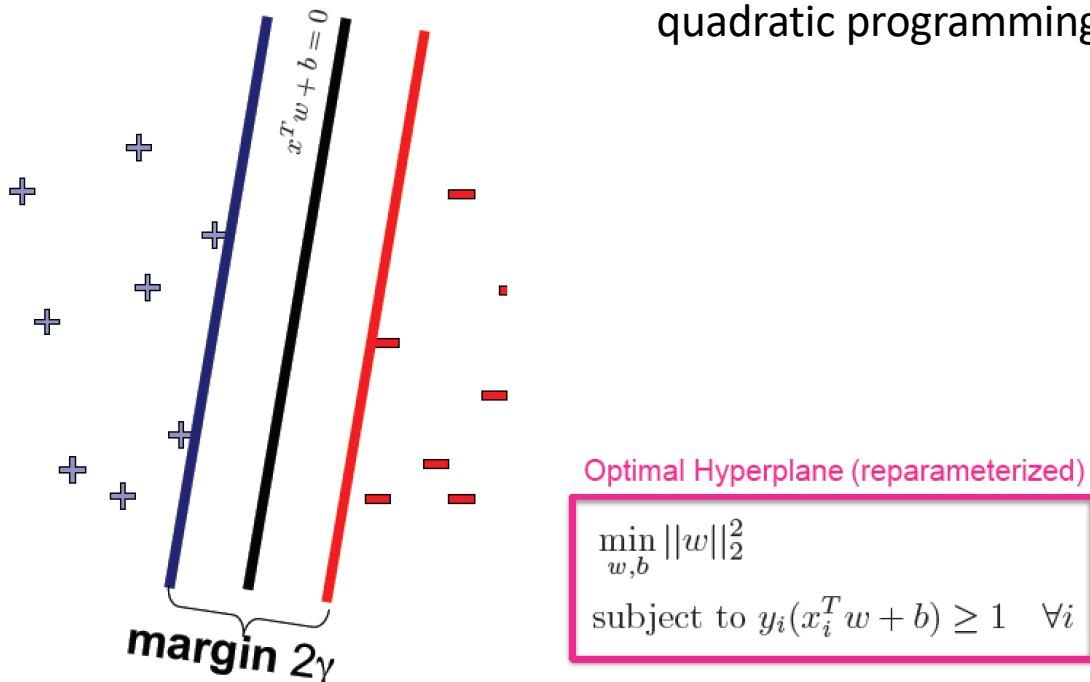


Our optimization problem

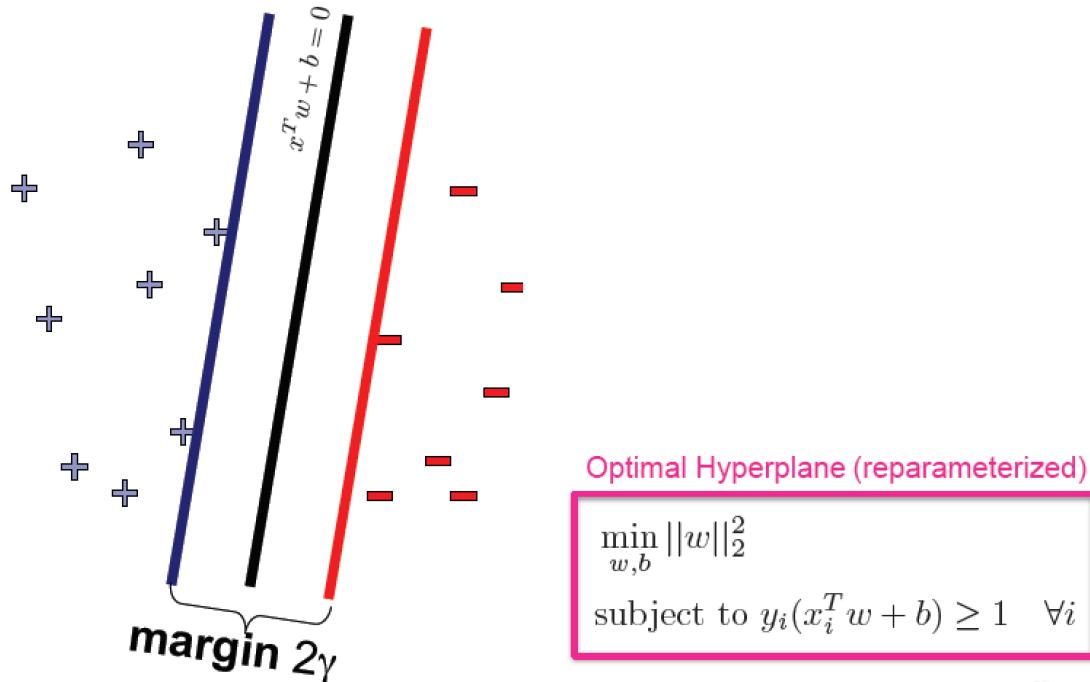


Final version

Solvable efficiently –
quadratic programming problem



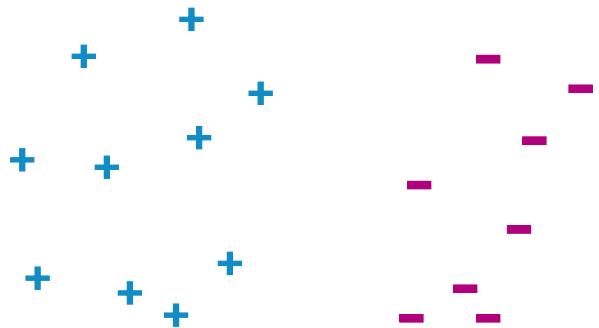
What are support vectors?



What if the data are not linearly separable?

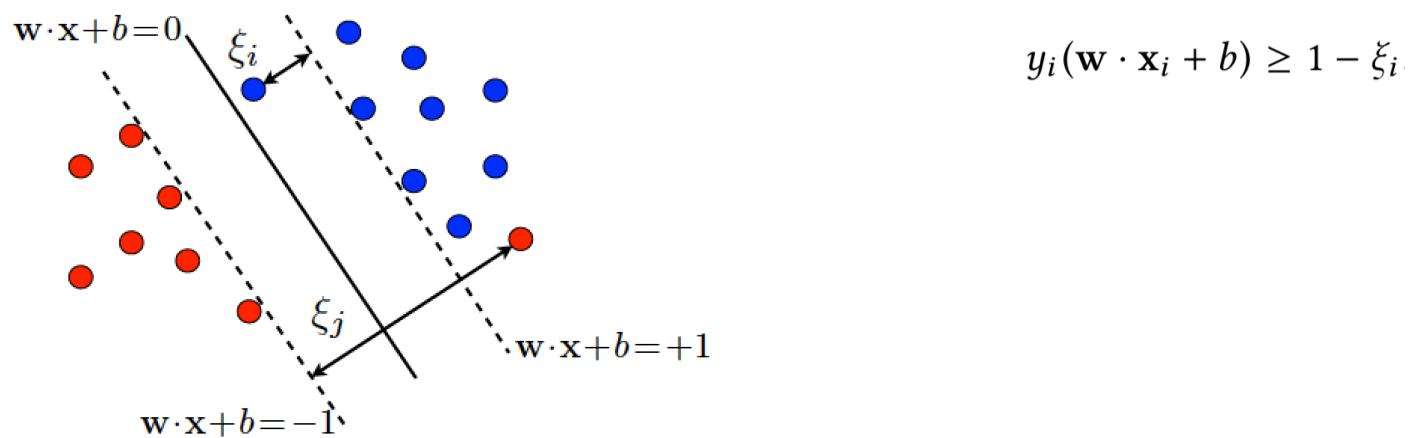
What if data are not linearly separable?

Use feature maps...



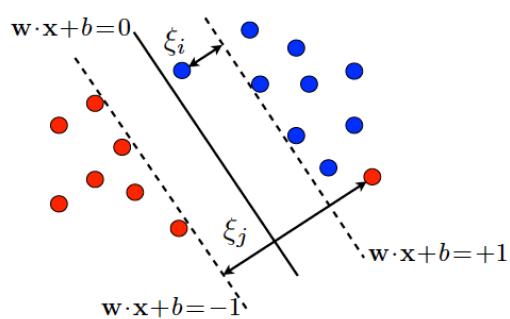
SVMs can be kernelized!!!!

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**

batch size of one (or online learning)

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

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:

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

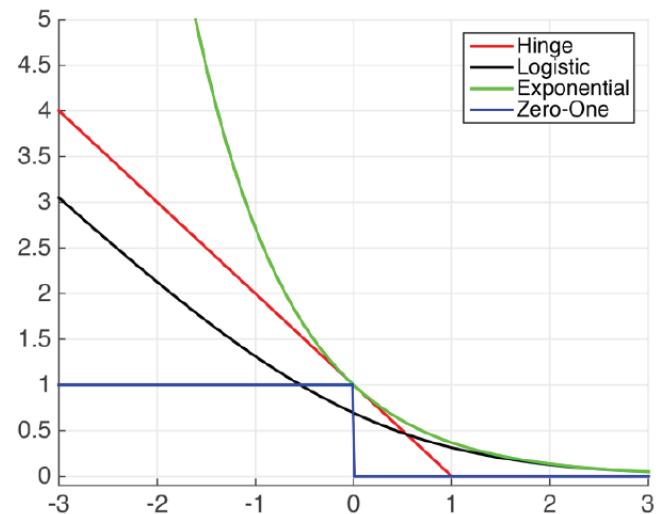
Logistic loss:

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

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