

Kernels

Machine Learning – CSE446

Carlos Guestrin

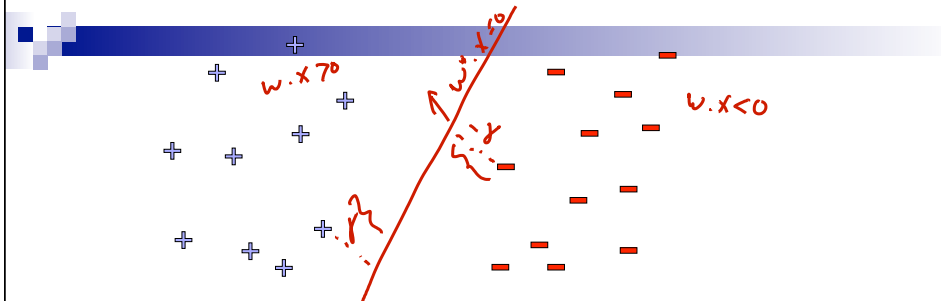
University of Washington

October 28, 2013

©Carlos Guestrin 2005-2013

1

Linear Separability: More formally, Using Margin



■ Data linearly separable, if there exists

□ a vector $\exists w^*$, $\|w^*\|=1$

□ a margin $\delta > 0$

■ Such that all points are at least δ away from $w \cdot x = 0$

$\forall t : f(y^{(t)} = +1, w^* \cdot x^{(t)} \geq \delta$

$y^{(t)} = -1, w^* \cdot x^{(t)} \leq -\delta$

Linearly separable:
 $y^{(t)} w^* \cdot x^{(t)} \geq \delta$

©Carlos Guestrin 2005-2013

2

Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:
 - Given a sequence of labeled examples: $(x^{(1)}, y^{(1)}) \dots (x^{(T)}, y^{(T)})$
not iid
 - Each feature vector has bounded norm:
 $\forall x^{(i)} \quad \|x^{(i)}\| \leq R$
 - If dataset is linearly separable:
 $\exists w^*, \|w^*\| = 1, \forall x^{(i)} \quad w^* \cdot x^{(i)} \geq \gamma$ for $\gamma > 0$
- Then the number of mistakes made by the online perceptron on any such sequence is bounded by

$\left(\frac{R}{\gamma}\right)^2$ wow!!

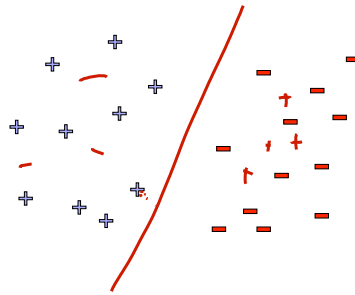
↖ constant, doesn't grow with T !!
😊

©Carlos Guestrin 2005-2013

3

Beyond Linearly Separable Case

- Perceptron algorithm is super cool!
 - No assumption about data distribution!
 - Could be generated by an oblivious 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 not linearly separable
 - Can't expect never to make mistakes again
 - Analysis extends to non-linearly separable case
 - Very similar bound, see Freund & Schapire
 - Converges, but ultimately may not give good accuracy (make many many many mistakes)



o: if data is very very non-linearly separable

©Carlos Guestrin 2005-2013

4

What if the data is not linearly separable?

Use features of features of features of features....

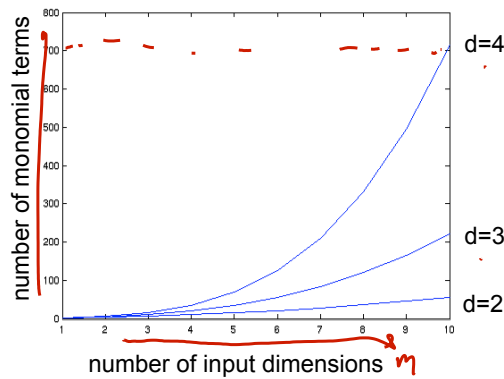
$$\Phi(\mathbf{x}) : \mathbb{R}^m \mapsto F$$

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ \dots \\ \sin x_1 \log x_2 \end{pmatrix}$$

Feature space can get really large really quickly!

Higher order polynomials

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!}$$



m – input features
d – degree of polynomial

Remark:
e.g. fit polynomials of any degree \approx same cost

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

Perceptron Revisited

- Given weight vector $w^{(t)}$, predict point x by:

$$\hat{y} = \text{sign}(w^{(t)} \cdot x)$$

- Mistake at time t : $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$

- Thus, write weight vector in terms of mistaken data points only: $w^{(t)} = 0$ *for simplicity*

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

$$w^{(t)} = \sum_{j \in M^{(t)}} y^{(j)} x^{(j)}$$

- Prediction rule now:

$$\text{sign}(x \cdot w^{(t)}) = \text{sign}\left(x \cdot \sum_{j \in M^{(t)}} y^{(j)} x^{(j)}\right) = \text{sign}\left(\sum_{j \in M^{(t)}} y^{(j)} x \cdot x^{(j)}\right)$$

only depends on dot prod. between x & mistakes

- When using high dimensional features:

$$\text{sign}(\phi(x) \cdot w^{(t)}) = \text{sign}\left(\sum_{j \in M^{(t)}} y^{(j)} \phi(x) \cdot \phi(x^{(j)})\right)$$

can often often be comp efficiently, even in high dim

©Carlos Guestrin 2005-2013

Dot-product of polynomials

$$u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

$$v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$\Phi(u) \cdot \Phi(v)$ ~~is~~ ^{for} polynomials of degree exactly d

$$d=1 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = u_1 v_1 + u_2 v_2 = u \cdot v$$

$$d=2 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix} \cdot \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix} = u_1^2 v_1^2 + 2u_1 u_2 v_1 v_2 + u_2^2 v_2^2 = (u_1 v_1 + u_2 v_2)^2 = (u \cdot v)^2$$

Proof by single of induction

For poly. of degree exactly d :

$$\phi(u) \cdot \phi(v) = (u \cdot v)^d = K(u, v)$$

Kernel trick!!

©Carlos Guestrin 2005-2013

8

Finally the Kernel Trick!!! (Kernelized Perceptron)

- Every time you make a mistake, remember $(\mathbf{x}^{(t)}, y^{(t)})$
↳ keep index $M(t)$ of all mistakes up to time t

- Kernelized Perceptron prediction for \mathbf{x} :

$$\underbrace{\text{sign}(\mathbf{w}^{(t)} \cdot \phi(\mathbf{x}))}_g = \sum_{j \in M^{(t)}} y^{(j)} \phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x})$$

$$= \sum_{j \in M^{(t)}} y^{(j)} k(\mathbf{x}^{(j)}, \mathbf{x})$$

remember mistakes & use kernels for dot products

©Carlos Guestrin 2005-2013

9

Polynomial kernels

- All monomials of degree d in $O(d)$ operations:
 $\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d =$ polynomials of degree exactly d

- How about all monomials of degree up to d ?

□ Solution 0: $\phi(u) \cdot \phi(v) = \sum_{i=0}^d \binom{d}{i} (u \cdot v)^i$

- Better solution:

$$(u \cdot v)^1 + (u \cdot v)^2 + (u \cdot v)^0 = (u \cdot v + 1)^2$$

proof by "induction"

$$\phi(u) \cdot \phi(v) = k(u, v) = (u \cdot v + 1)^d$$

wow!!!

©Carlos Guestrin 2005-2013

10

Common kernels

MANY KERNELS ARE EFFICIENT

- Polynomials of degree exactly d

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

- Polynomials of degree up to d

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

- Gaussian (squared exponential) kernel ✓

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

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

Radial basis function
RBF kernel
Gaussian kernel
↳ infinite dim feature space

©Carlos Guestrin 2005-2013

11

What you need to know

- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proofs
- The kernel trick
- Kernelized Perceptron
- Derive polynomial kernel
- Common kernels
- In online learning, report averaged weights at the end

©Carlos Guestrin 2005-2013

12

Your Midterm...

- Content: Everything up to last Wednesday (Perceptron)...
- Only 80mins, so arrive early and settle down quickly, we'll start and end on time
- “Open book”
 - Textbook, Books, Course notes, Personal notes, your HWs
- Bring a calculator that can do log ☺
- No:
 - Computers, tablets, phones, other materials, internet devices, wireless telepathy or wandering eyes...
- The exam:
 - Covers key concepts and ideas, work on understanding the big picture, and differences between methods

©Carlos Guestrin 2005-2013

13

Support Vector Machines

Machine Learning – CSE446

Carlos Guestrin

University of Washington

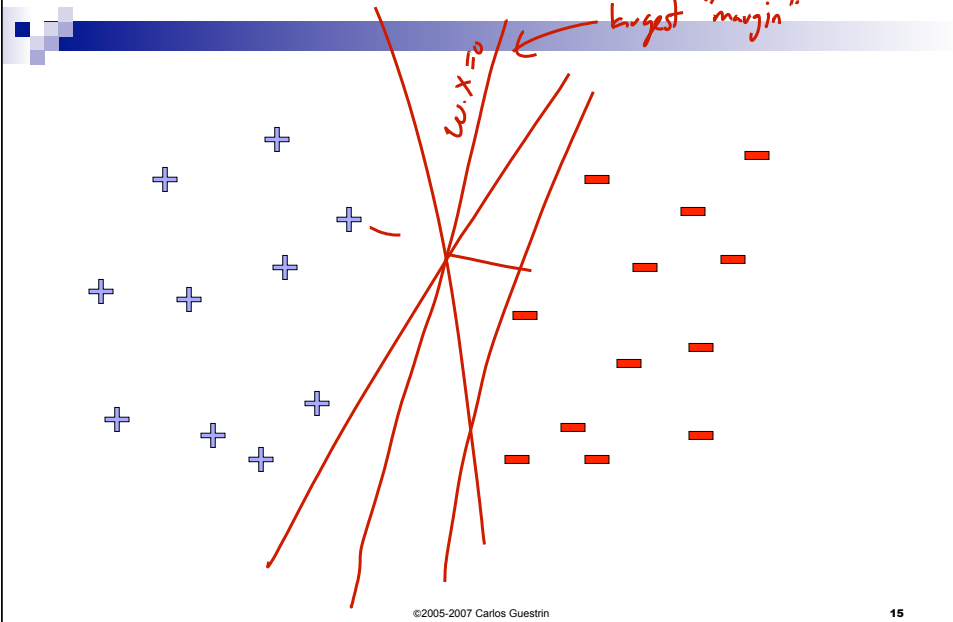
October 28, 2013

©Carlos Guestrin 2005-2013

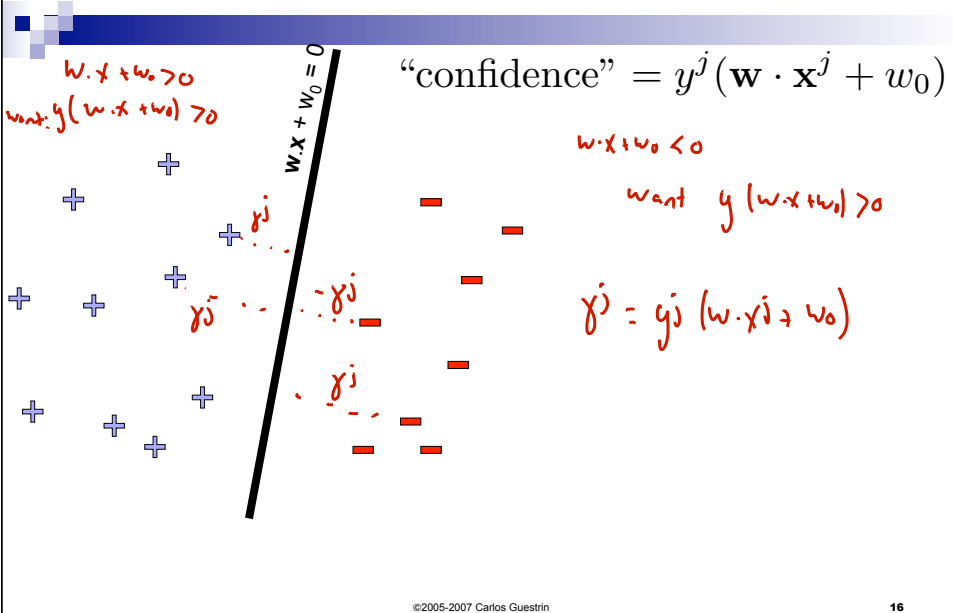
14

\approx perception + L2 regularization.

Linear classifiers – Which line is better?

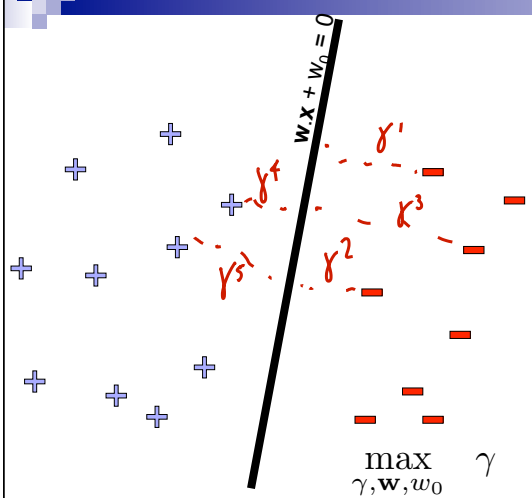


Pick the one with the largest margin!



Maximize the margin

SVN: $\max_{w, w_0} \min_j y_j$



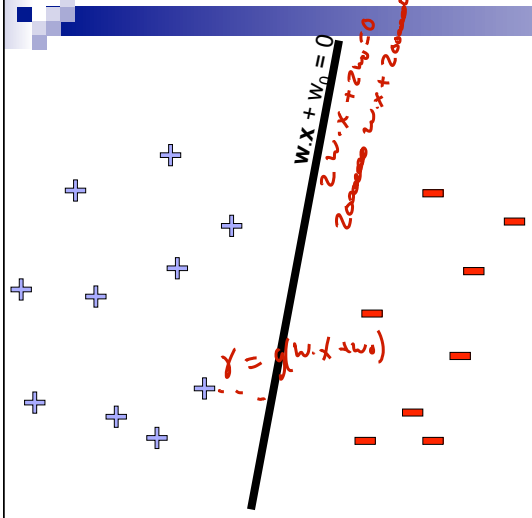
$\max_{w, w_0} \min_j y_j$
 $y_j = y_j (w \cdot x^j + w_0)$

find me w, w_0 such that all points have margin at least γ

$\max_{\gamma, w, w_0} \gamma$

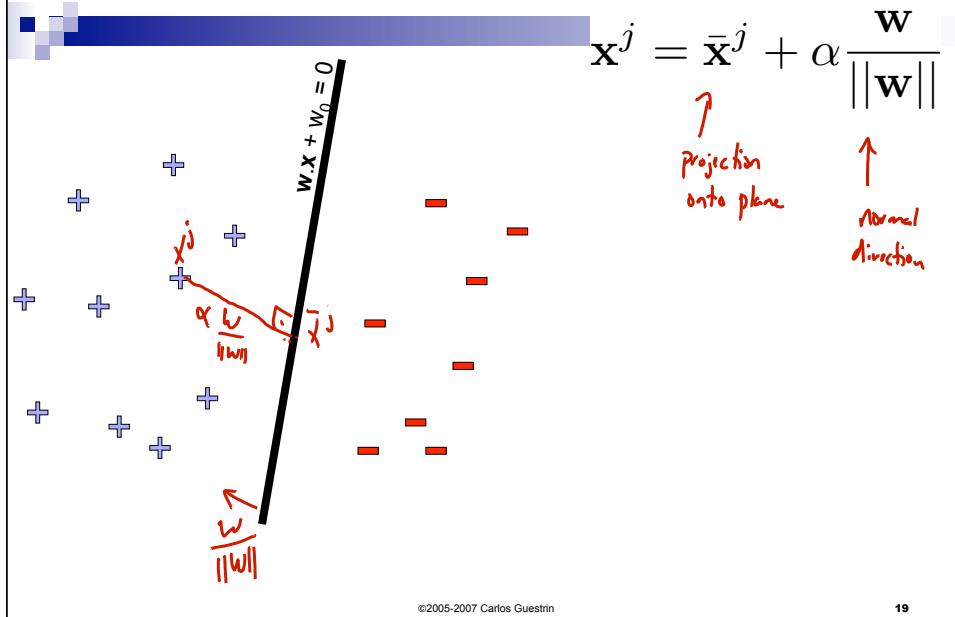
$y^j (w \cdot x^j + w_0) \geq \gamma, \forall j \in \{1, \dots, N\}$

But there are many planes...

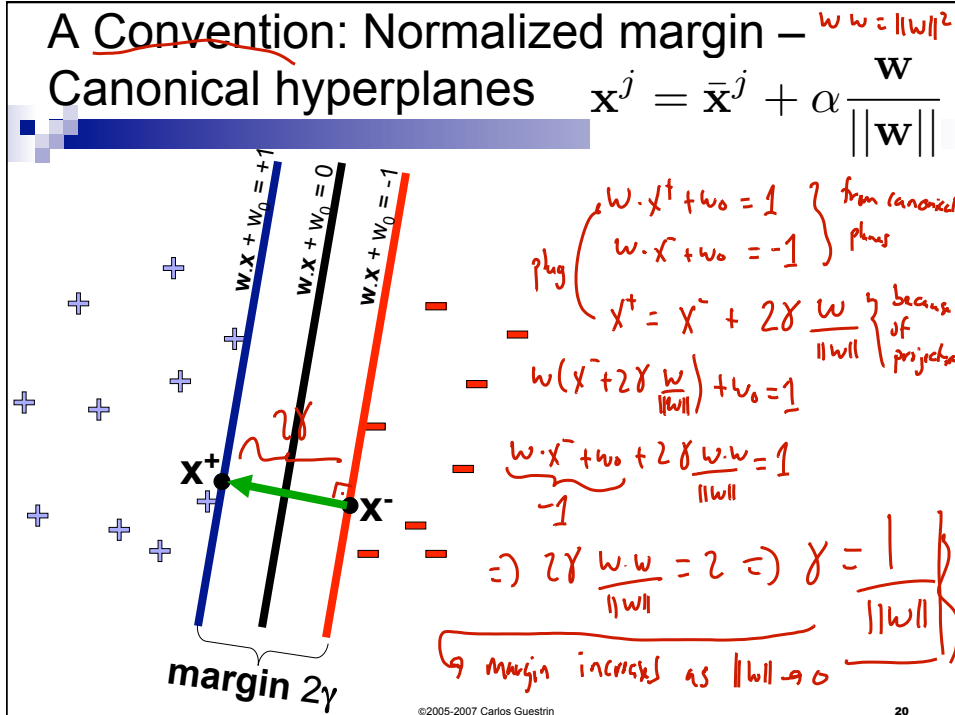


γ will increase in magnitude as we increase w, w_0
 \Rightarrow unbounded problem

Review: Normal to a plane



A Convention: Normalized margin – Canonical hyperplanes



Margin maximization using canonical hyperplanes

$\max \frac{1}{f} \equiv \min f$
 $f > 0 \implies \min f^2$

Unnormalized problem: $\max_{\gamma, \mathbf{w}, w_0} \gamma$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq \gamma, \forall j \in \{1, \dots, N\}$

Normalized Problem:
 $\max \gamma \equiv \max \frac{1}{\|\mathbf{w}\|} \equiv \min \|\mathbf{w}\| \equiv \min \|\mathbf{w}\|^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1$

$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j \in \{1, \dots, N\}$

©2005-2007 Carlos Guestrin 21

Support vector machines (SVMs)

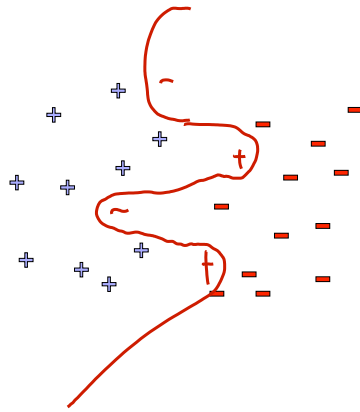
$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$
 $y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j \in \{1, \dots, N\}$

- Solve efficiently by many methods, e.g.,
 - quadratic programming (QP)
 - Well-studied solution algorithms
 - Stochastic gradient descent
- Hyperplane defined by support vectors
 - ↗ could forget all data except for support vectors & get same solution

Don't change solution

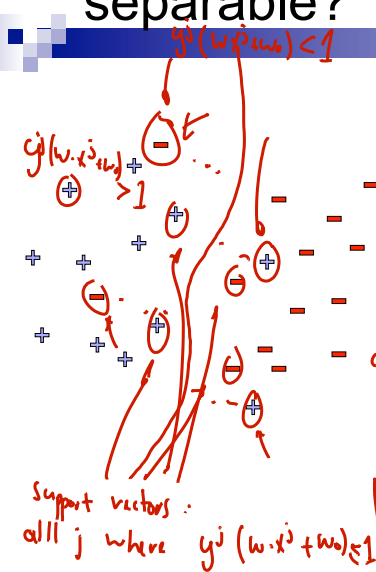
©2005-2007 Carlos Guestrin 22

What if the data is not linearly separable?



Use features of features of features of features....

What if the data is still not linearly separable?



$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

$$y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \geq 1, \forall j$$

- If data is not linearly separable, some points don't satisfy margin constraint:

$$\exists j \quad y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) < 1 \Leftrightarrow 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) > 0$$

- How bad is the violation?

$$\text{constraint violation} = \begin{cases} 0 & 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) \leq 0 \\ 1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0) & \text{otherwise} \end{cases} \left\{ \begin{array}{l} \text{a violation of constraint} \\ (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+ \end{array} \right.$$

- Tradeoff margin violation with $\|\mathbf{w}\|$:

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

↑
C > 0 slack penalty

SVMs for Non-Linearly Separable meet my friend the Perceptron...

- Perceptron was minimizing the hinge loss:

$$\sum_{j=1}^N (-y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

hinge loss

- SVMs minimize the regularized hinge loss!!

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

regularization

Stochastic Gradient Descent for SVMs

- Perceptron minimization:

$$\sum_{j=1}^N (-y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SGD for Perceptron:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \mathbb{1}_{[y^{(t)}(\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)} \leq 0)]} y^{(t)} \mathbf{x}^{(t)}$$

update weights

$\eta=1$

mistake
 $\nabla (-y^{(t)}(\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)} + w_0))_+$

- SVMs minimization:

$$\|\mathbf{w}\|_2^2 + C \sum_{j=1}^N (1 - y^j (\mathbf{w} \cdot \mathbf{x}^j + w_0))_+$$

- SGD for SVMs:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)}$$

$$+ \eta \left(C \mathbb{1}_{[1 - y^{(t)}(\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)} + w_0) \leq 0]} y^{(t)} \mathbf{x}^{(t)} - 2 \mathbf{w}^{(t)} \right)$$

step size is important in SVMs

pick C? by cross validation

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
 - Hinge loss
 - A.K.A. adding slack variables
- SVMs = Perceptron + L2 regularization
- Can also use kernels with SVMs
- Can optimize SVMs with SGD
 - Many other approaches possible