

Support vector machines (SVMs)

Lecture 5

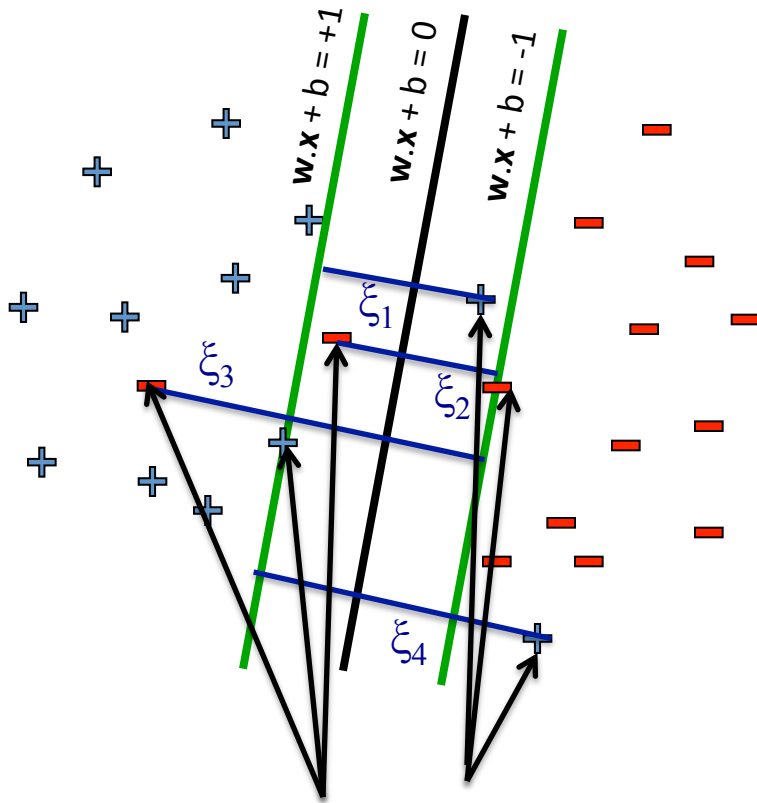
David Sontag
New York University

Soft margin SVM

$$\operatorname{argmin}_{\mathbf{w}, \xi_i \geq 0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$$

“slack variables”



Slack penalty $C > 0$:

- $C = \infty \rightarrow$ minimizes upper bound on 0-1 loss
- $C \approx 0 \rightarrow$ points with $\xi_i = 0$ have big margin
- **Select using cross-validation**

Support vectors:

Data points for which the constraints are binding

Soft margin SVM

QP form:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{w}, \xi_i \geq 0} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i \end{aligned}$$

More “natural” form:

$$\operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$$

where:

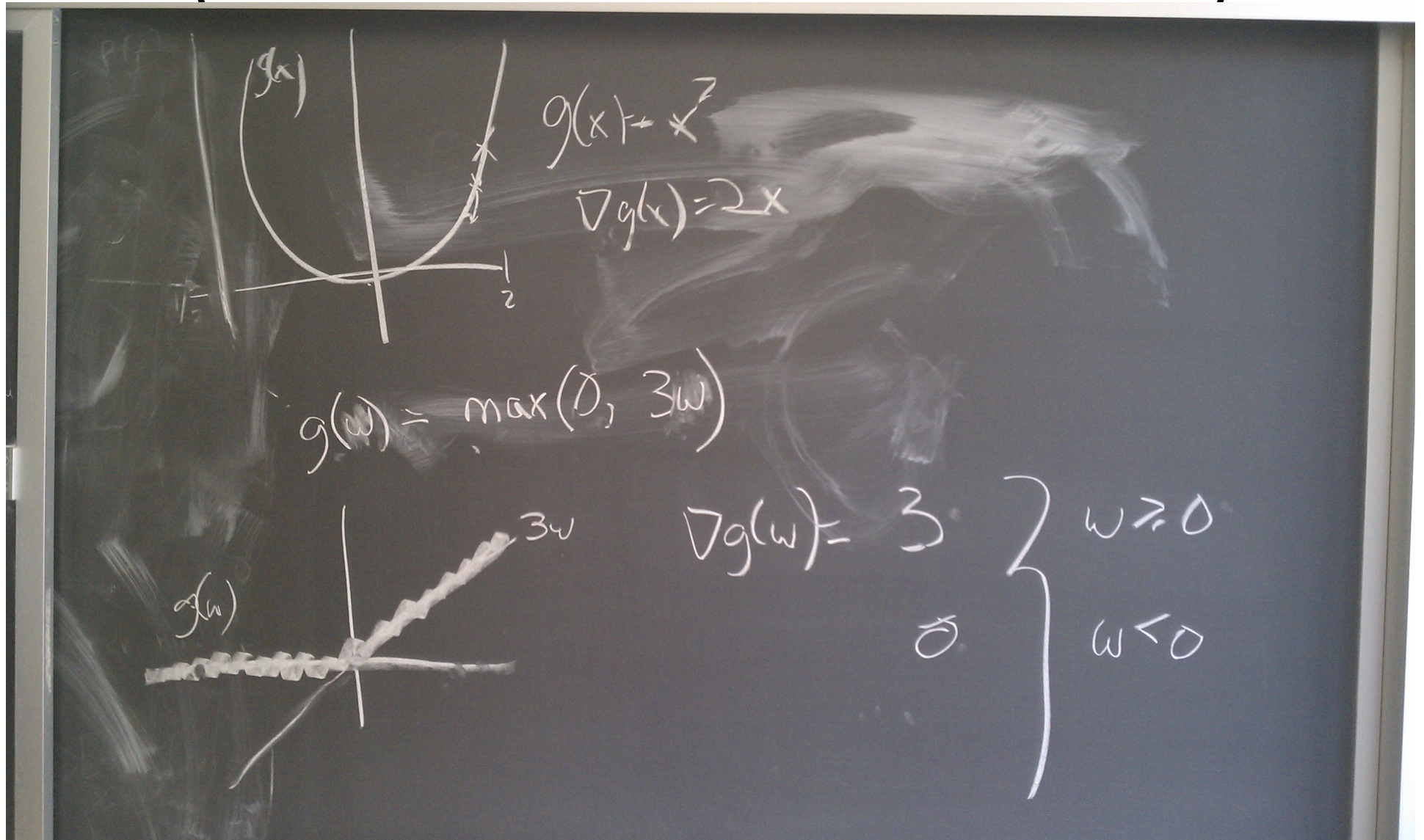
Equivalent if
 $C = \frac{1}{m\lambda}$

$$f(\mathbf{w}) \stackrel{\text{def}}{=} \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{Regularization term}} + \underbrace{\frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\}}_{\text{Empirical loss}}$$

Regularization
term

Empirical loss

Subgradient (for non-differentiable functions)



(Sub)gradient descent of SVM objective

$$f(\mathbf{w}) \stackrel{\text{def}}{=} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\}$$

Handwritten notes on a chalkboard showing the subgradient of the SVM objective function and the update rule for the weight vector.

$$\nabla f(\omega) = \lambda \vec{\omega} - \frac{1}{m} \sum_{\substack{i \text{ s.t.} \\ y_i \langle \omega, \mathbf{x}_i \rangle < 1}} y_i \frac{\vec{1}}{\mathbf{x}_i}$$

$\vec{\omega}^0 = \vec{0}$
for $t=1, \dots$

$$\vec{\omega}^{t+1} = \vec{\omega}^t - \eta^t \nabla f(\omega^t)$$

Step size:

$$\eta_t = \frac{1}{t\lambda}$$

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0$, $t=0$

While not converged

$t = t+1$

 Choose a stepsize, η_t

 Choose a direction, p_t

 Go!

 Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0$, $t=0$

For iter = 1,2,...,20

For $j=1,2,\dots, |data|$

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t x_j) < 1$

$w_{t+1} = (1-\eta_t\lambda) w_t + \eta_t y_j x_j$

Else

$w_{t+1} = (1-\eta_t\lambda) w_t$

Output: w_{t+1}

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0, t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0, t=0$

For iter = 1,2,...,20

For j=1,2,..., |data|

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t \cdot x_j) < 1$

$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$

Else

$w_{t+1} = w_t - \eta_t \lambda w_t$

Output: w_{t+1}

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0, t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0, t=0$

For iter = 1,2,...,20

For j=1,2,..., |data|

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t x_j) < 1$

$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$

Else

$w_{t+1} = w_t - \eta_t \lambda w_t$

Output: w_{t+1}

Convergence choice : Fixed number of iterations

$T=20 * |data|$

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0$, $t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0$, $t=0$

For iter = 1,2,...,20

For $j=1,2,\dots, |data|$

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t \cdot x_j) < 1$

$$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$$

Else

$$w_{t+1} = w_t - \eta_t \lambda w_t$$

Output: w_{t+1}

Stepsize choice: - Initialize with $1/\lambda$
- Decays with $1/t$

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0$, $t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0$, $t=0$

For iter = 1,2,...,20

For $j=1,2,\dots, |data|$

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t x_j) < 1$

$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$

Else

$w_{t+1} = w_t - \eta_t \lambda w_t$

Output: w_{t+1}

Direction choice: Stochastic approx to the subgradient

Subgradient calculation

Objective: $\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_i \max\{0, 1 - y_i w \cdot x_i\}$

Stochastic Approx: $\frac{\lambda}{2} \|w\|^2 + \max\{0, 1 - y_i w \cdot x_i\}$

For a randomly chosen data point i

*(in the assignment the choice of i is **not random** - easier to debug and compare between students).*

Subgradient calculation

Objective: $\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_i \max\{0, 1 - y_i w \cdot x_i\}$

Stochastic Approx: $\frac{\lambda}{2} \|w\|^2 + \max\{0, 1 - y_i w \cdot x_i\}$

(sub)gradient:

$$\lambda \|w\| + \frac{d}{dw} \max\{0, 1 - y_i w \cdot x_i\}$$

Subgradient calculation

Objective:

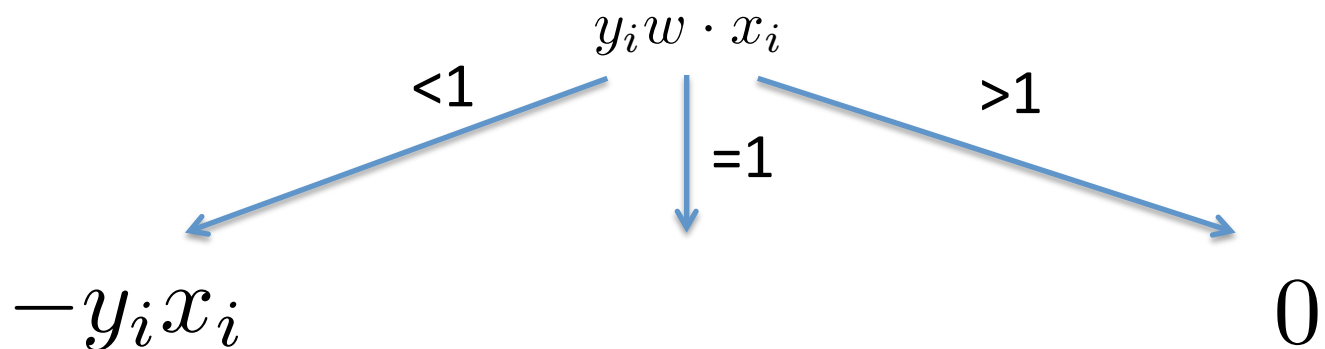
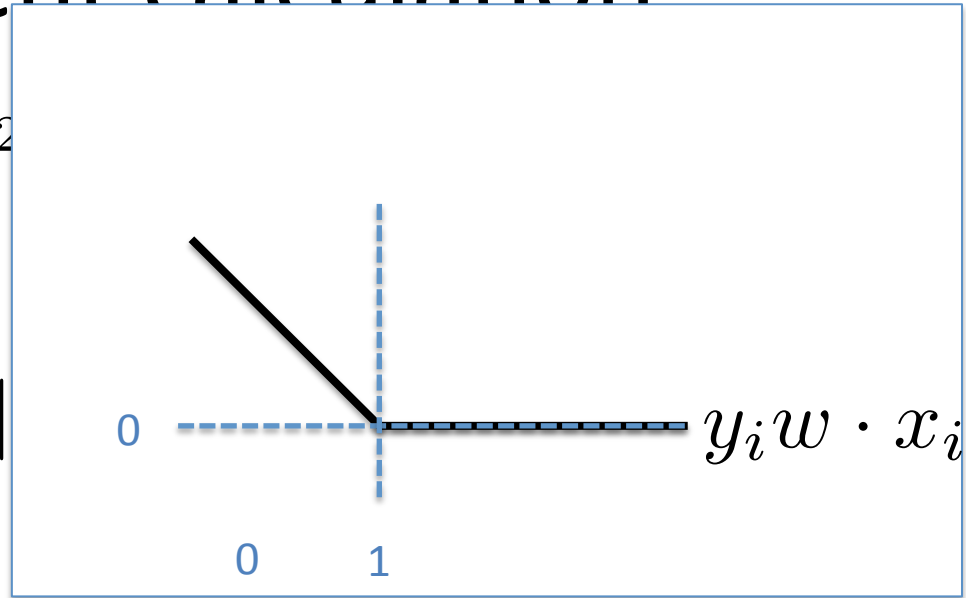
$$\frac{\lambda}{2} \|w\|^2$$

Stochastic Approx:

$$\frac{\lambda}{2} \|w\|^2$$

(sub)gradient:

$$\lambda \|w\| + \frac{d}{dw} \max\{0, 1 - y_i w \cdot x_i\}$$



Subgradient calculation

Objective:

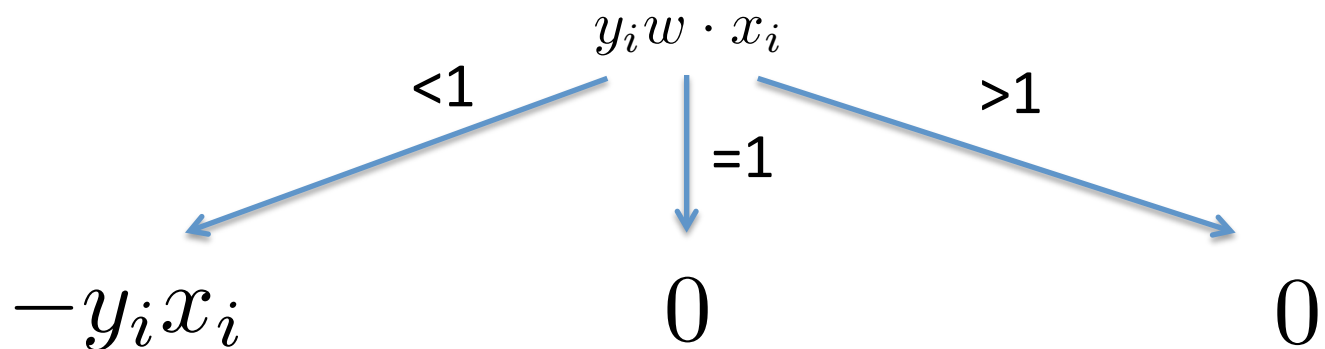
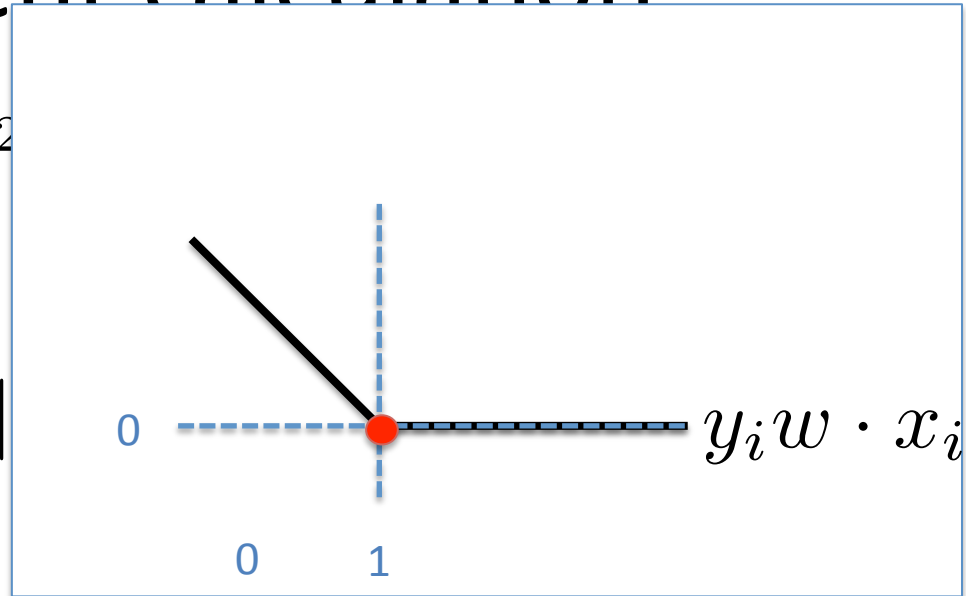
$$\frac{\lambda}{2} \|w\|^2$$

Stochastic Approx:

$$\frac{\lambda}{2} \|w\|^2$$

(sub)gradient:

$$\lambda \|w\| + \frac{d}{dw} \max\{0, 1 - y_i w \cdot x_i\}$$



Subgradient calculation

Objective: $\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_i \max\{0, 1 - y_i w \cdot x_i\}$

Stochastic Approx: $\frac{\lambda}{2} \|w\|^2 + \max\{0, 1 - y_i w \cdot x_i\}$

(sub)gradient:

if $y_i w \cdot x_i < 1$ $\lambda w - y_i x_i$

else $\lambda w + 0$

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0, t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0, t=0$

For iter = 1,2,...,20

For $j=1,2,\dots, |\text{data}|$

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t \cdot x_j) < 1$

$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$

Else

$w_{t+1} = w_t - \eta_t(\lambda w_t + 0)$

Output: w_{t+1}

Direction choice: Stochastic approx to the subgradient

if $y_i w \cdot x_i < 1$

$\lambda w - y_i x_i$

else

$\lambda w + 0$

The Pegasos Algorithm

General framework

Initialize: $w_1 = 0, t=0$

While not converged

$t = t+1$

Choose a stepsize, η_t

Choose a direction, p_t

Go!

Test for convergence

Output: w_{t+1}

Pegasos Algorithm (from homework)

Initialize: $w_1 = 0, t=0$

For iter = 1,2,...,20

For $j=1,2,\dots, |data|$

$t = t+1$

$\eta_t = 1/(t\lambda)$

If $y_j(w_t x_j) < 1$

$w_{t+1} = w_t - \eta_t(\lambda w_t - y_j x_j)$

Else

$w_{t+1} = w_t - \eta_t \lambda w_t$

Output: w_{t+1}

Go: update $w_{t+1} = w_t - \eta_t p_t$

Why is this algorithm interesting?

- Simple to implement, state of the art results.
 - Notice similarity to Perceptron algorithm!
Algorithmic differences: updates if insufficient margin, scales weight vector, and has a learning rate.
- Since based on *stochastic* gradient descent, its running time guarantees are probabilistic.
- Highlights interesting tradeoffs between running time and data.

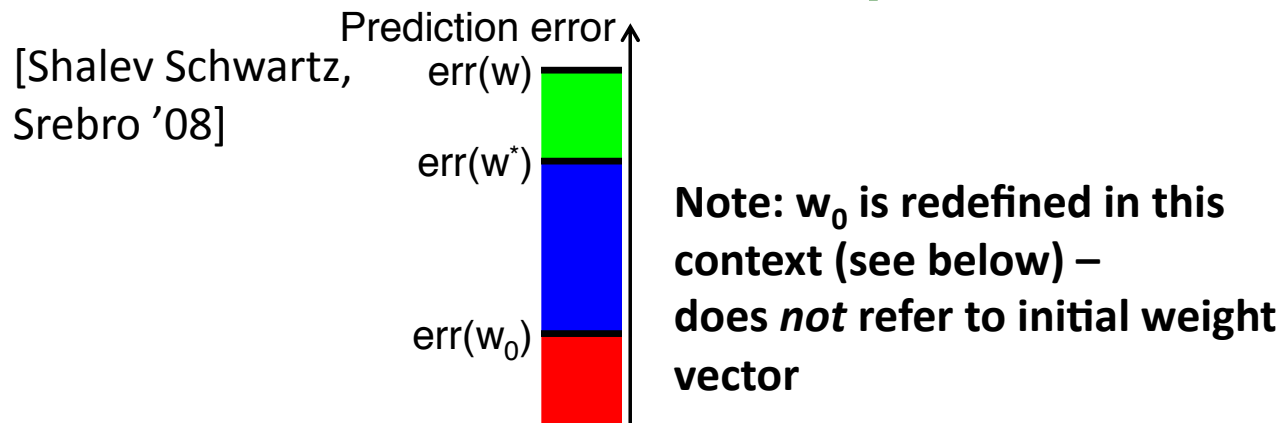
Much faster than previous methods

- **3 datasets** (provided by Joachims)
 - Reuters CCAT (800K examples, 47k features)
 - Physics ArXiv (62k examples, 100k features)
 - Covertypes (581k examples, 54 features)

Training Time
(in seconds):

	Pegasos	SVM-Perf	SVM-Light
Reuters	2	77	20,075
Covertypes	6	85	25,514
Astro-Physics	2	5	80

Approximate algorithms

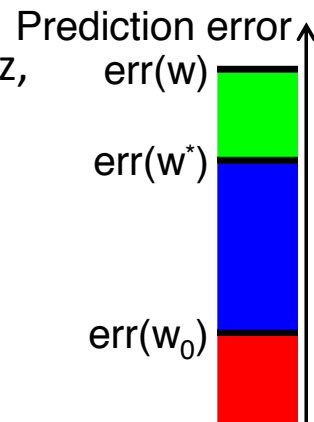


- **Approximation error:**
 - Best error achievable by large-margin predictor
 - Error of population minimizer
 $w_0 = \operatorname{argmin} E[f(w)] = \operatorname{argmin} \lambda \|w\|^2 + E_{x,y}[\operatorname{loss}(\langle w, x \rangle; y)]$
- **Estimation error:**
 - Extra error due to replacing $E[\operatorname{loss}]$ with empirical loss
 $w^* = \operatorname{argmin} f_n(w)$
- **Optimization error:**
 - Extra error due to only optimizing to within finite precision

From ICML'08 presentation (available [here](#))

Approximate algorithms

[Shalev Schwartz,
Srebro '08]



- **Approximation error:**
 - Best error achievable by large-margin
 - Error of population minimizer
 $w_0 = \operatorname{argmin} E[f(w)] = \operatorname{argmin} \lambda \|w\|^2$
- **Estimation error:**
 - Extra error due to replacing $E[\text{loss}]$ by $f_n(w)$
 $w^* = \operatorname{argmin} f_n(w)$
- **Optimization error:**
 - Extra error due to only optimizing to

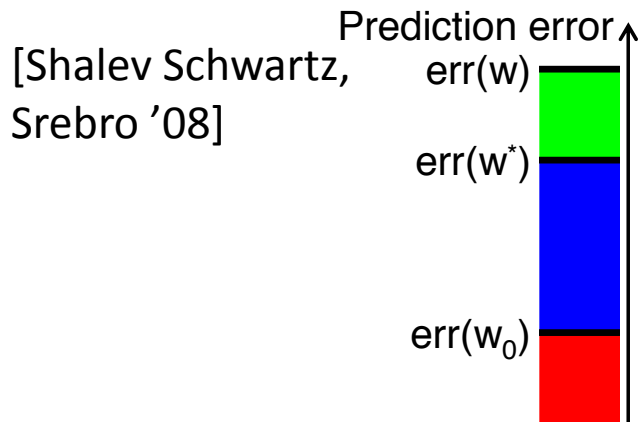
Pegasos Guarantees

After $T = \tilde{O}\left(\frac{1}{\delta \lambda \epsilon}\right)$ updates:

$$\operatorname{err}(w_T) < \operatorname{err}(w_0) + \epsilon$$

With probability $1 - \delta$

Approximate algorithms



Running time does **NOT** depend on:

-# training examples!

It **DOES** depend on:

- Dimensionality d (why?)
- Approximation ϵ and δ
- Difficulty of problem λ

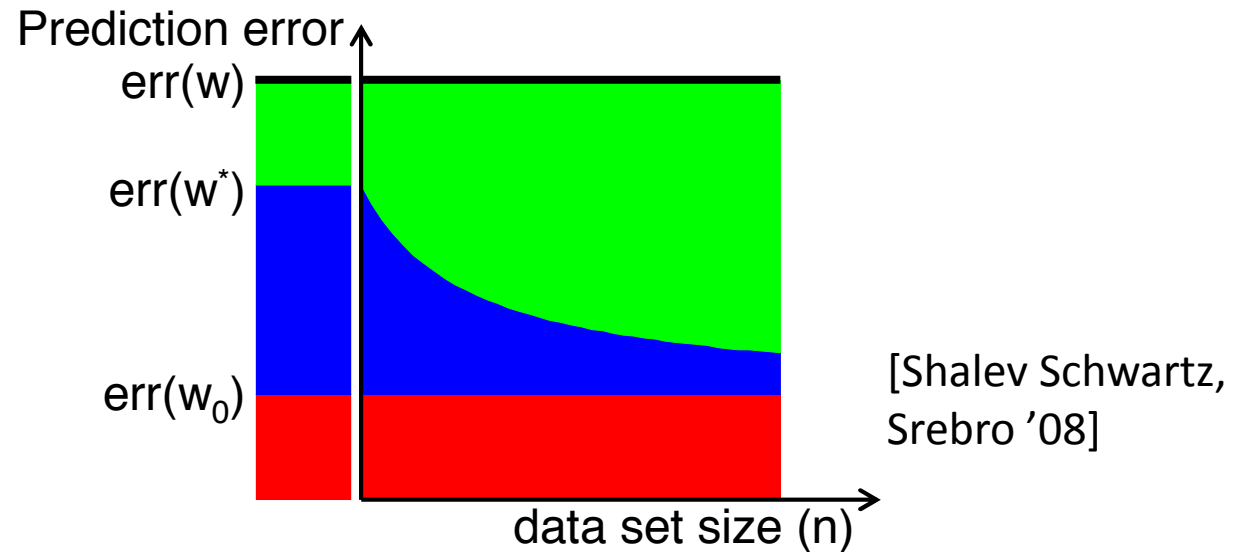
Pegasos Guarantees

After $T = \tilde{O}\left(\frac{1}{\delta\lambda\epsilon}\right)$ updates:

$$err(w_T) < err(w_0) + \epsilon$$

With probability $1 - \delta$

But how is that possible?



As the dataset grows, our approximations can be worse to get the same error!