# Lecture #3: Online Learning*

**Janardhan Rao (Jana) Doppa**

School of EECS, Washington State University

* Slides partly based on Koby Crammer

# Formal setting – Classification

- **Instances**

  $$\mathbf{x} \in \mathcal{X}$$

  - emails

- **Labels**

  $$y \in \mathcal{Y} = \{-1 \; ; \; 1\}$$

  - Spam vs. non-spam

- **Prediction rule**

  $$f(\mathbf{x}) = \widehat{y}$$

  - Linear prediction rule

- **Loss**

  $$\ell(\widehat{y}, y) \in \mathbb{R}_+$$

  - No. of mistakes

# Predictions

- Continuous predictions : $f : \mathcal{X} \to \mathbb{R}$

  - Label $\mathsf{sign}(f(\mathbf{x}))$

  - Confidence $|f(\mathbf{x})|$

- Linear Classifiers

  - Prediction :
  $$\begin{aligned} \hat{y} &= \mathsf{sign}(f(\mathbf{x})) \\ &= \arg\max_{y \in \mathcal{Y}} \mathbf{w} \cdot \Phi(\mathbf{x}, y) \\ &= \mathsf{sign}(\mathbf{w} \cdot \mathbf{x}) \end{aligned}$$
  $$|f(\mathbf{x})| = |\mathbf{w} \cdot \mathbf{x}|$$

3

# Loss Functions

- **Natural Loss:**
  - Zero-One loss

$$\ell(\widehat{y}, y) = \begin{cases} 0 & y = \widehat{y} \\ 1 & y \neq \widehat{y} \end{cases}$$
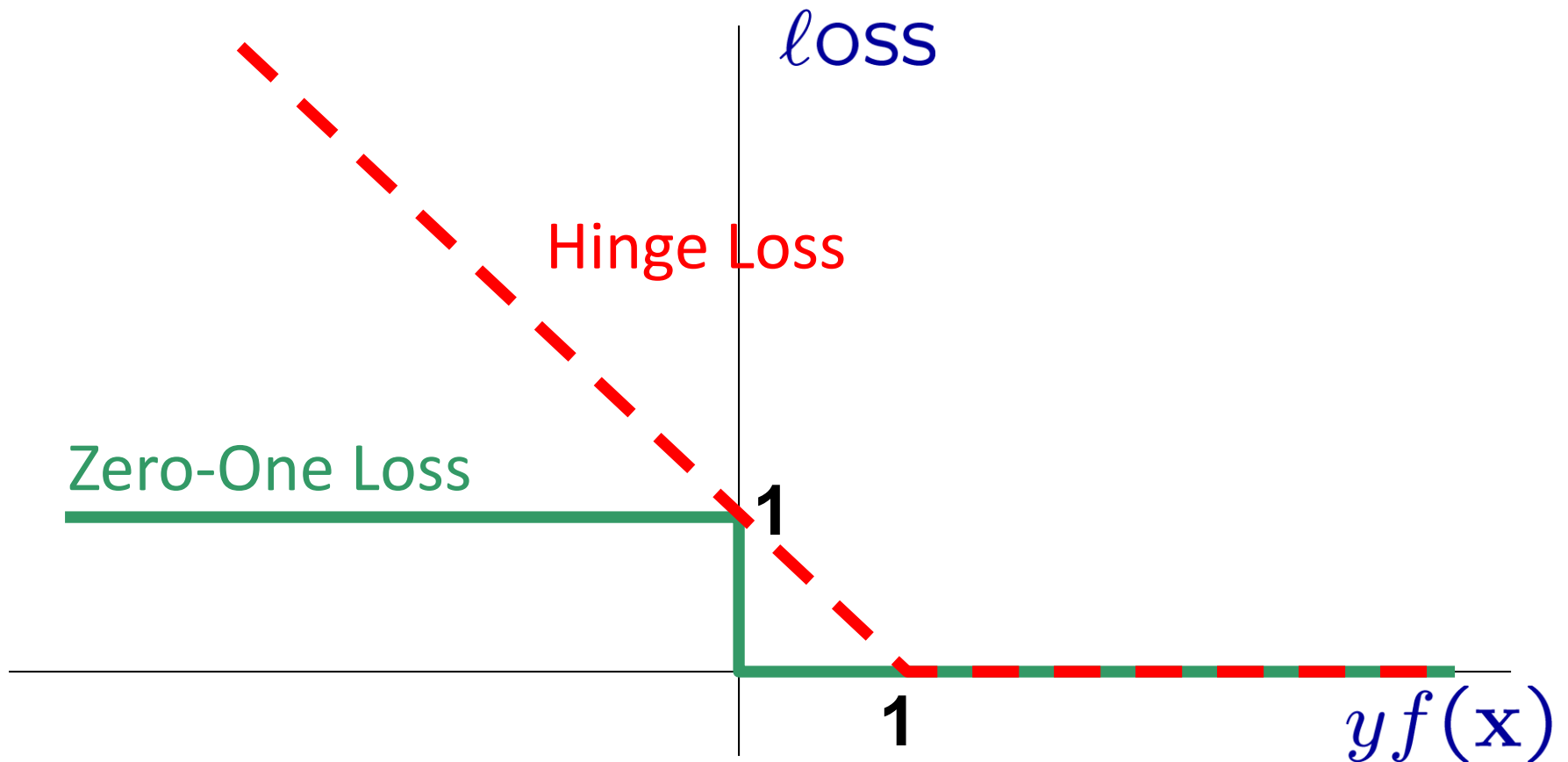
- **Real-valued-predictions loss:**
  - Hinge loss

$$\ell(\widehat{y}, y) = \max\{0, 1 - yf(\mathbf{x})\}$$

  - Exponential loss (Boosting)

# Loss Functions

# Online Framework

- Initialize Classifier $f_1(\mathbf{x})$

- Algorithm works in rounds $t = 1 \ldots T \ldots$

- On round $t$ the online algorithm :

  - Receives an input instance $\mathbf{x}_t$
  - Outputs a prediction $f_t(\mathbf{x}_t) = \widehat{y}_t$
  - Receives a feedback label $y_t$
  - Computes loss $\ell(\widehat{y}_t, y_t)$
  - Updates the prediction rule $f_t \rightarrow f_{t+1}$

- Goal :

  - Suffer small cumulative loss $\sum_t \ell(\widehat{y}_t, y_t)$

# Margin

- ***Margin*** of an example $(\mathbf{x}_t, y_t)$ with respect to the classifier $\mathbf{w}_t$ :
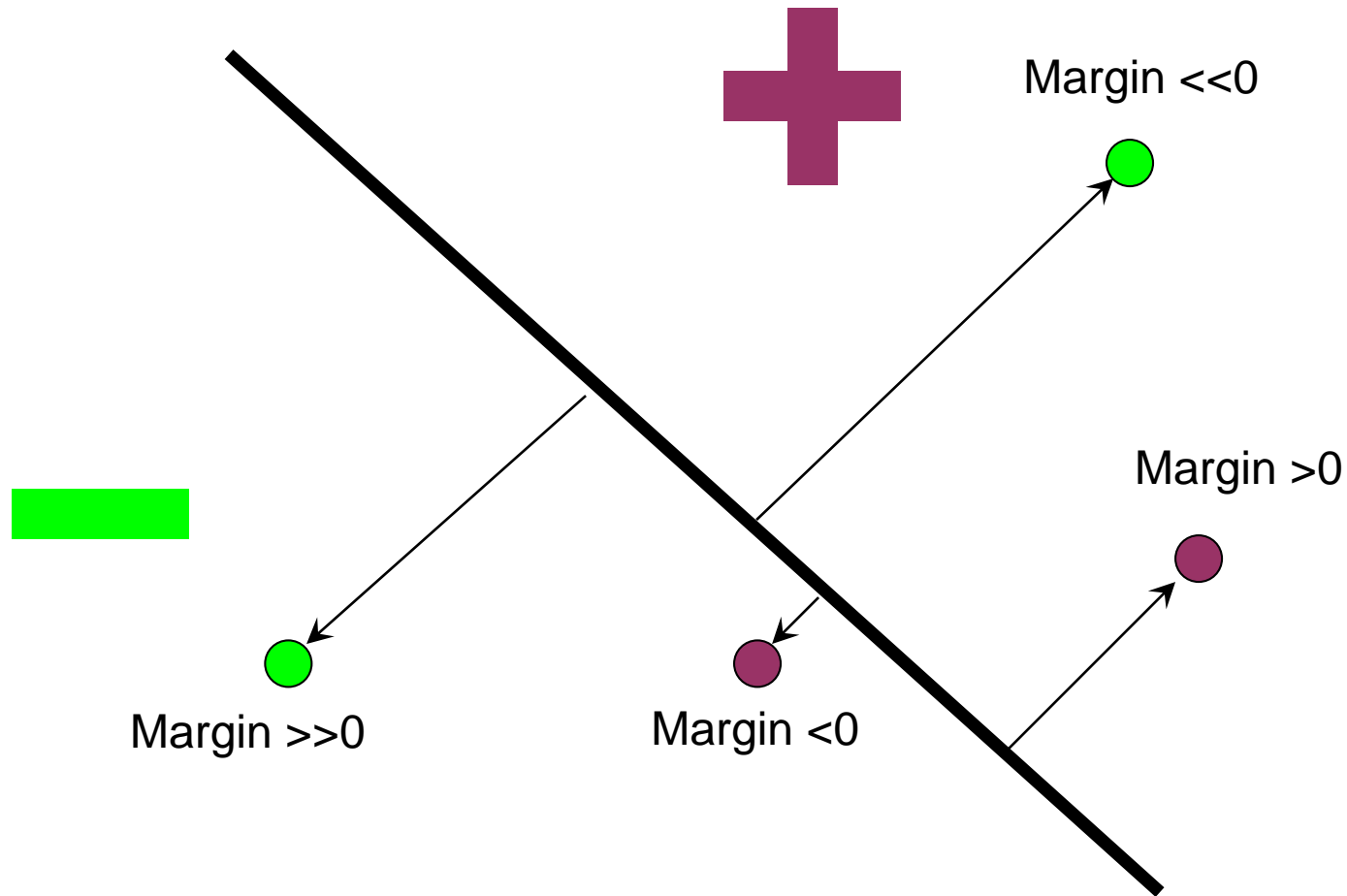
$$y_t(\mathbf{w}_t \cdot \mathbf{x}_t)$$

- Note :

$$y_t(\mathbf{w}_t \cdot \mathbf{x}_t) > 0 \quad \Leftrightarrow \quad \ell_{01}(\widehat{y}_t, \mathbf{x}_t) = 0$$
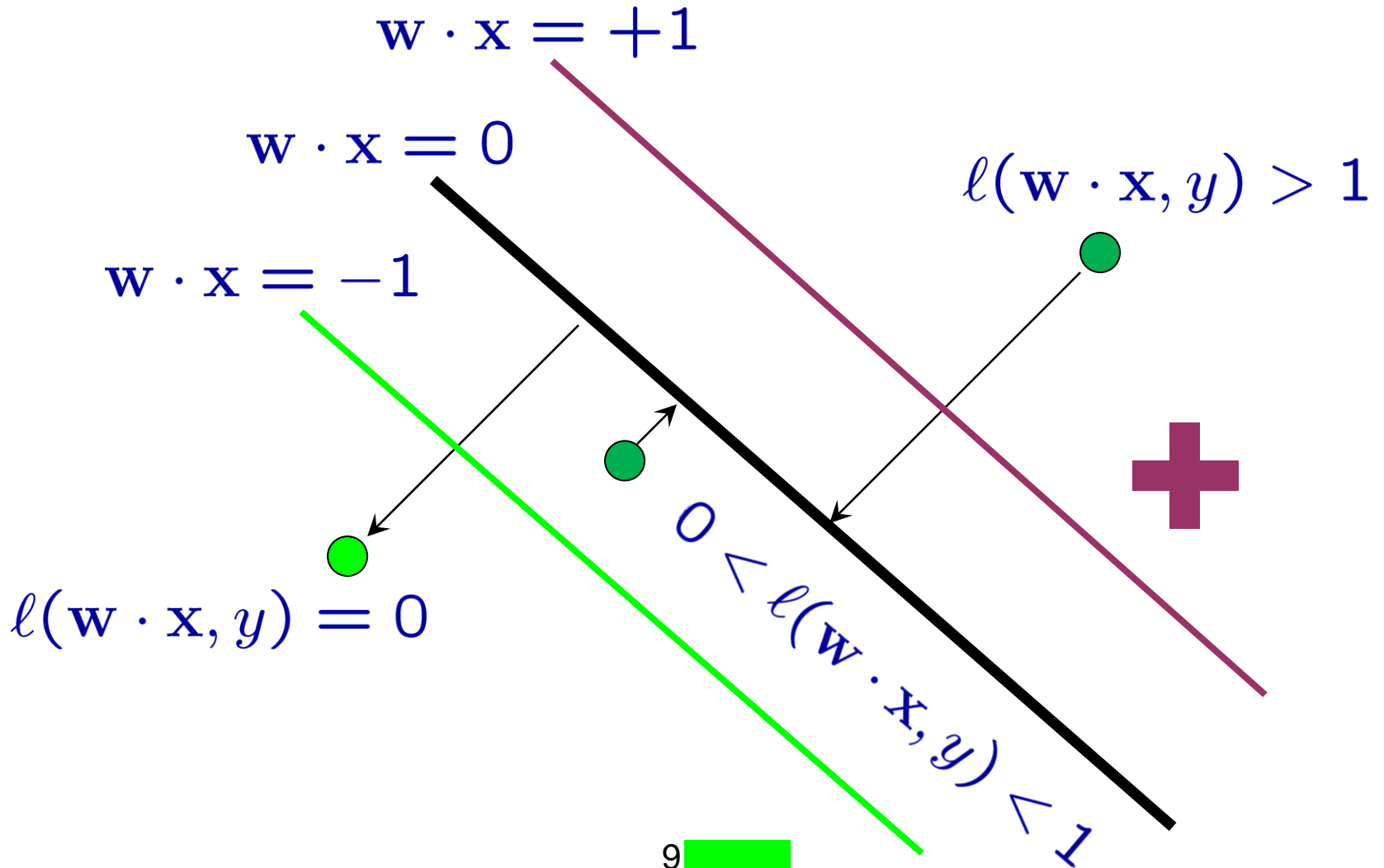
- The set $(\mathbf{x}_1, y_1) \ldots (\mathbf{x}_T, y_T)$ is separable iff there exists *u* such that

$$y_t(\mathbf{u} \cdot \mathbf{x}_t) > 0 \qquad \forall t$$

# Geometrical Interpretation

Margin <<0

Margin >0

Margin >>0

Margin <0

# Hinge Loss

$$\mathbf{w} \cdot \mathbf{x} = +1$$

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

$$\ell(\mathbf{w} \cdot \mathbf{x}, y) > 1$$

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

$$\ell(\mathbf{w} \cdot \mathbf{x}, y) = 0$$

$$0 < \ell(\mathbf{w} \cdot \mathbf{x}, y) < 1$$

9

# **Why Online Learning?**

- Fast
- Memory efficient - process one example at a time
- Simple to implement
- Formal guarantees – Mistake bounds
- Online to Batch conversions
- No statistical assumptions
- Adaptive

# **Update Rules**

- Online algorithms are based on an update rule which defines $f_{t+1}$ from $f_t$ (and possibly other information)

- **Linear Classifiers** : find $\mathbf{w}_{t+1}$ from $\mathbf{w}_t$ based on the input $(\mathbf{x}_t, y_t)$

- **Some Update Rules** :

  - Perceptron (Rosenblat)
  - ALMA (Gentile)
  - ROMMA (Li & Long)
  - NORMA (Kivinen et. al)

  - MIRA (Crammer & Singer)
  - EG (Littlestown and Warmuth)
  - Bregman Based (Warmuth)
  - CWL (Dredze et. al)

11

# Design Principles of Algorithms

- If the learner suffers non-zero loss at any round, then we want to balance two goals:

  - **Corrective:** Change weights so that we don't make this error again

  - **Conservative:** Don't change the weights too much

$$d(w_{t+1}, w_t) + \eta L(y_t, w_{tt})$$

# The Perceptron Algorithm

- If the learner suffers non-zero loss at any round, then we want to balance two goals:

$$d(w_{t+1}, w_t) + \eta L(y_t, w_{tt})$$

- **Euclidean distance**

- **Hinge loss**

- Stochastic Gradient Descent (SGD)

# The Perceptron Algorithm ($\eta = 1$)

- If No-Mistake $\quad y_t(\mathbf{w}_t \cdot \mathbf{x}_t) > 0$

  - Do nothing $\quad \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$

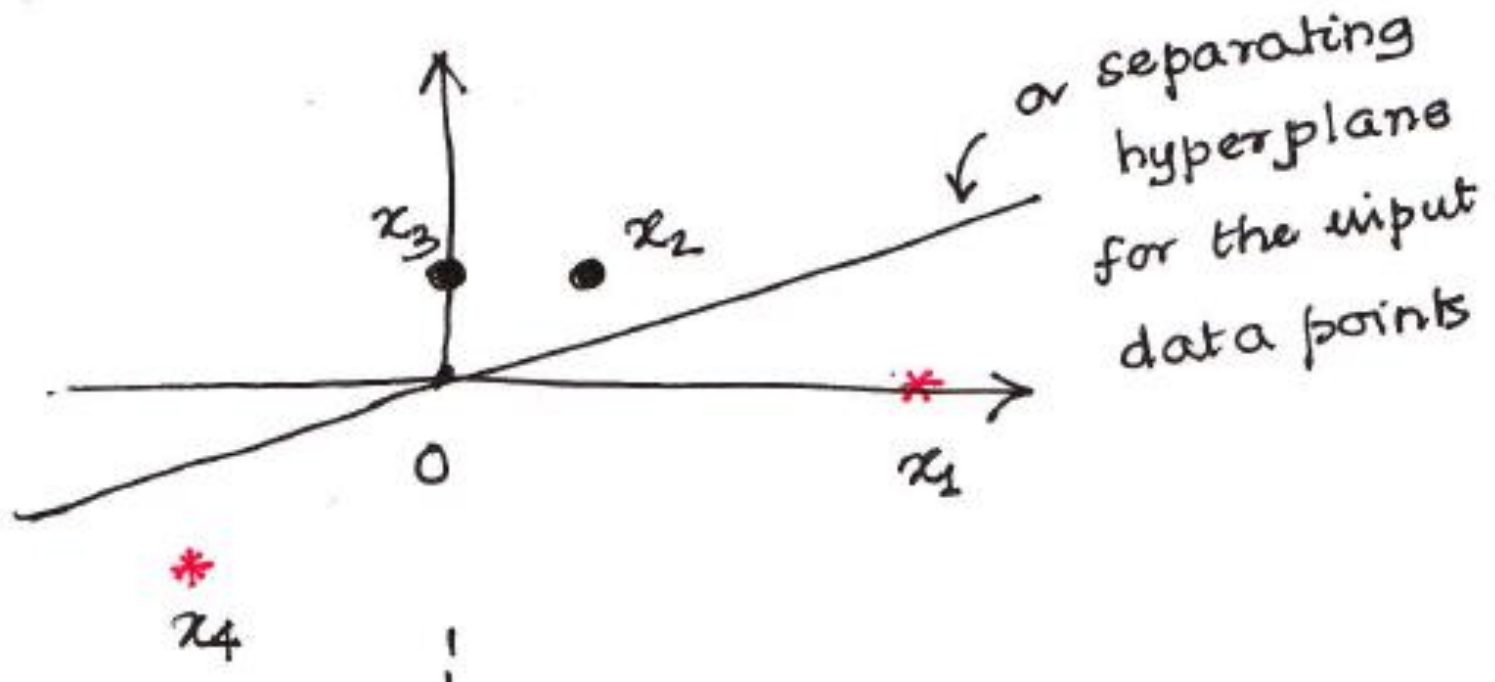- If Mistake $\quad y_t(\mathbf{w}_t \cdot \mathbf{x}_t) \leq 0$

  - Update $\quad \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{x}_t$

# The Perceptron Algorithm ($\eta = 1$)

- When mistake happens, what does the update do?

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

  ✓ $w_{t+1}$ moves "closer to" $x_t$ OR

  ✓ $x_t$ moves towards the positive side of the decision boundary

  ○ **If $y_t = -1$:**  $$w_{t+1} = w_t - x_t$$

  ✓ $w_{t+1}$ moves "away from" $x_t$ OR

  ✓ $x_t$ moves towards the negative side of the decision boundary

- In both cases, we are moving towards the "correct solution"

Training Data:  $((4,0), 1)$,  $((1,1), -1)$,  $((0,1), -1)$, $((-2,-2), 1)$



a separating hyperplane for the input data points

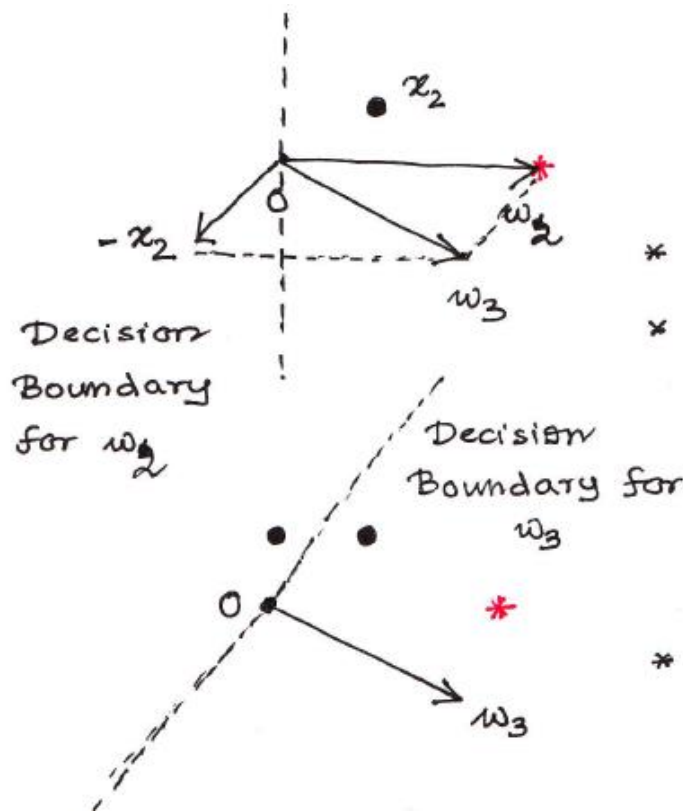# Running Example #1

Round 1:
* $w_1 = 0$

Training Data: $((4,0), 1), ((1,1), -1), ((0,1), -1), ((-2,-2), 1)$

* $y_t \langle w_t, x_t \rangle$ for $t = 1$
  $= 0$ as well.

* $w_2 = w_1 + y_1 x_1$
  $= (4, 0)$



Decision Boundary for $w_2$

Decision Boundary for $w_3$

Round 2:

* $y_2 \langle w_2, x_2 \rangle < 0$

* So $w_3 = w_2 + y_2 x_2$
  $= (4, 0) - (1, 1) = (3, -1)$
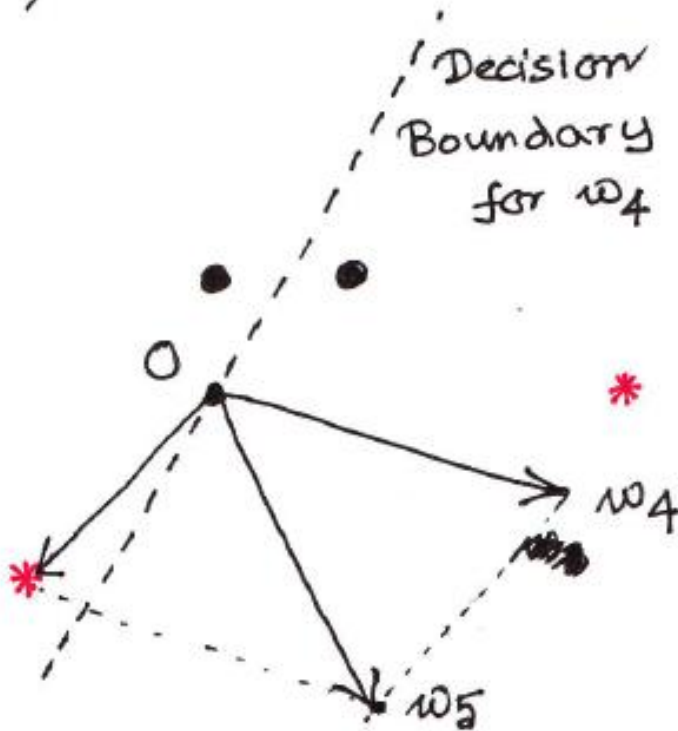
Round 3:

* $y_3 \langle w_3, x_3 \rangle > 0$
  Correct. $w_4 = w_3 = (3, -1)$

# Running Example #1

Training Data: $((4,0), 1), ((1,1), -1), ((0,1), -1), ((-2,-2), 1)$



Decision
Boundary
for $w_4$

Round 4:

* $y_4 \langle w_4, x_4 \rangle < 0$

* So $w_5 = w_4 + y_4 x_4$
$$= (3, -1) + (-2, -2)$$
$$= (1, -3)$$

# Running Example #2

$( (1, 1), 1 ), \quad ( (1, -1), -1 ), \quad ( (-1, 1), -1 ), \quad ( (-1, -1), 1 )$



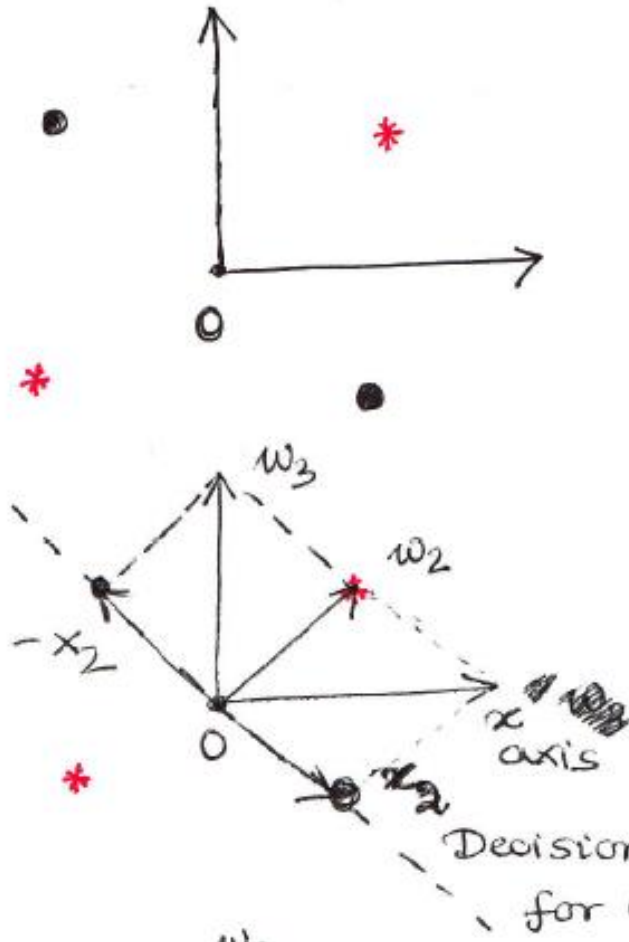Initially, $w_1 = 0$. In round 1,

* $y_1 \langle w_1, x_1 \rangle \leq 0$, so mistake.

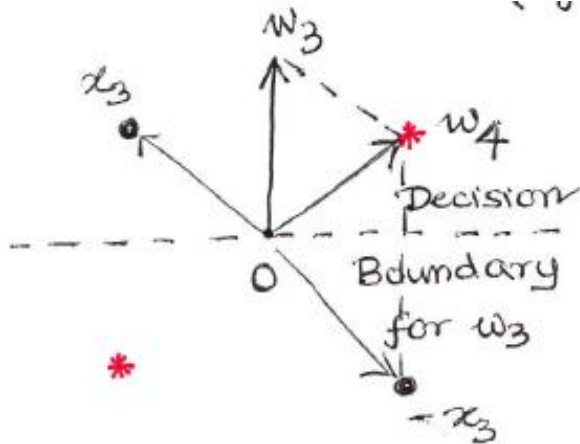* $w_2 = w_1 + y_1 x_1 = (1, 1)$

Round 2:

* $y_2 \langle w_2, x_2 \rangle \leq 0$

* $w_3 = w_2 + y_2 x_2 = (1, 1) - (1, -1)$

$= (0, 2)$
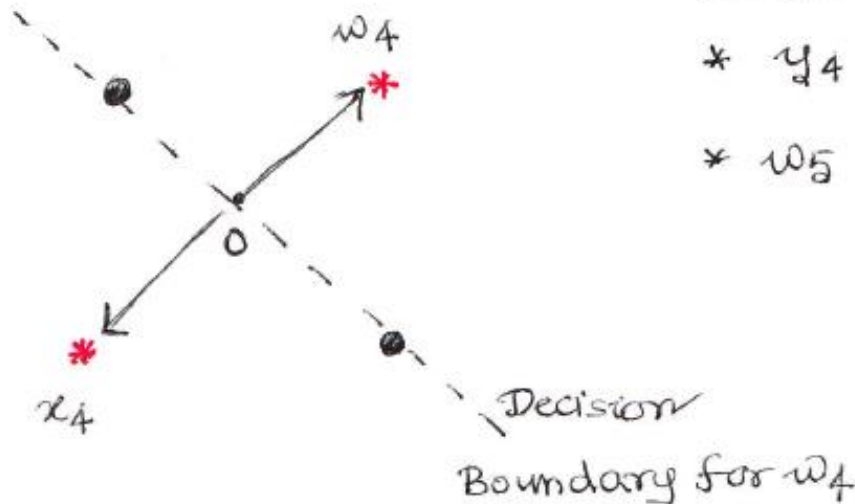
19

# Running Example #2



Round 3:

* $y_3 \langle w_3, x_3 \rangle \leq 0$

* $w_4 = w_3 + y_3 x_3 = (0, 2) - (-1, 1)$

  $= (1, 1)$

Round 4:

* $y_4 \langle w_4, x_4 \rangle \leq 0.$

* $w_5 = w_4 + y_4 x_4 = (1, 1) + (-1, -1)$

  $= 0$

# The Perceptron Algorithm

- **Suppose $w_t$ makes a mistake on $(x_t, y_t)$, and we update $w_{t+1}$ as $w_{t+1} = w_t + y_t x_t$. Is it possible for $w_{t+1}$ to also make a mistake on $(x_t, y_t)$ ?**

# The Perceptron Algorithm ($\eta = 1$)

- **Suppose $w_t$ makes a mistake on $(x_t, y_t)$, and we update $w_{t+1}$ as $w_{t+1} = w_t + y_t x_t$. Is it possible for $w_{t+1}$ to also make a mistake on $(x_t, y_t)$ ?**
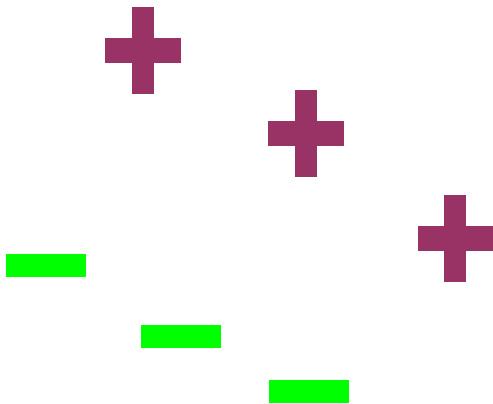
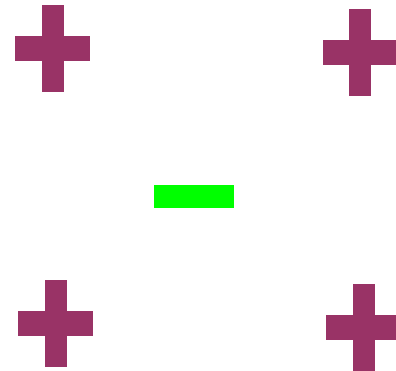  - ▲

  - ▲ Yes, depends on the learning rate $\eta$

# When does Perceptron converge?

- **Linear Separability**
  - There exists a hyper-plane (weight vector) separating the positive and negative points

Linearly separable

Not linearly separable

# Measure of Separability

- **Margin**
  - For a weight vector $w$, and training set $S$, margin of $w$ with respect to $S$ is defined as follows:

$$\gamma(w) = min_{(x,y) \in S} \quad y(w.x)$$

- The training data $S$ is linearly separable if there exists *at least* one weight vector $w$ for which the margin is positive, i.e., $\gamma(w) > 0$.

25

# Margin: Examples

**Low margin data**

**High margin data**

# **Perceptron: Convergence Result**

- **<u>Theorem:</u>** If the training data is linearly separable with margin $\gamma$, and if $\|x_i\| \le 1$ for all examples $(x_i, y_i)$ in the training set, then perceptron makes $\le \frac{1}{\gamma^2}$ mistakes.
  - ▲ Proof??

- Lower margin implies more mistakes

- May need more than one pass over the training data to get a classifier with no mistakes

# What if data is not linearly separable?

- **Ideally, we want to find a linear separator that makes the minimum number of mistakes on the training data**
  - NP-Hard problem! (Minsky and Papert, 1969)
  - This result killed the neural networks research in 1970's

- **Perceptron still works**
  - there will be few mistakes close to the decision boundary
  - will never converge to a single $w$ as we make more passes

# Problems with Perceptron

- Doesn't converge with inseparable data

- Weight updates may often be very "bold"

- Doesn't optimize margin

- Sensitive to the order of examples

  - **Voted and Averaged perceptron**

# Voted Perceptron

- **Initialization:** $m = 1; \; w_1 = 0; \; c_m = 1$

- **Training Examples:** for $t = 1, 2, 3, \ldots$
  - If mistake, update weights
    - $w_{m+1} = w_m + y_t x_t$
    - $m = m + 1$
    - $c_m = 1$
  - Else
    - $c_m = c_m + 1$  // counting how long $w_m$ survived

- **Output:** $(w_1, c_1), (w_2, c_2), (w_3, c_3), \ldots$

# Voted Perceptron Classifier

$$f(x) = sign\left(\sum_{i=1}^{m} c_i \, sign(<w_i, x>)\right)$$

- Any drawbacks of voted perceptron?

# Voted Perceptron Classifier

$$f(x) = sign\left(\sum_{i=1}^{m} c_i \, sign(< w_i, x >)\right)$$

- Any drawbacks of voted perceptron?

- Yes, we have to store all the classifiers (in practice could be many)

- How can we solve this problem?

33

# Averaged Perceptron

- Same algorithm as voted perceptron, but the classification rule is different

$$f_{average}(x) = sign\left(\sum_{i=1}^{m}(<c_i w_i, x>)\right)$$

$$f_{voted}(x) = sign\left(\sum_{i=1}^{m} c_i \, sign(<w_i, x>)\right)$$

# Averaged vs. Voted Perceptron

- Simple Example: If $c_1 = c_2 = c_3 = 1$

$$f_{average}(x) = sign(\langle w_1 + w_2 + w_3, x \rangle)$$

$$f_{voted}(x) = majority\ sign\ of\ \langle w_1, x \rangle, \langle w_2, x \rangle, \langle w_3, x \rangle$$

# **Some Practical Tricks**

- **Shuffling**
  - shuffling the training examples in each iteration

- **Variable learning rate**
  - decrease as learning progresses
  - follow some schedule
  - Set automatically by line search (converges faster)
  - See Leon Bottou's SGD website:
    http://leon.bottou.org/projects/sgd

- **Averaged Perceptron can be implemented very efficiently** (See Algorithm 7 in Hal's chapter)

# Perceptron vs. Averaged Perceptron

- Perceptron
  - $w_{t+1} = w_t + \eta . y_t . x_t$
  - $X_t$ has "n" features, but only "d" features are non-zero => $X_t$ is sparse
  - $O(n)$ => $O(d)$ If d << n, then it is a huge saving

- Averaged Perceptron
  - $w = 0$, $w_{sum} = 0$, count $= 0$
  - If we make a mistake, then we do the following:
    - Update $w = w + \delta$
    - Update $w_{sum} = w_{sum} + w$ => $O(n)$
    - Update count = count + 1
  - In the end, $w_{avg} = w_{sum} / count$

# Some Practical Tricks

- **Learning Curve**
  - Training iterations vs. number of mistakes
  - You want to see that the mistakes decrease as we increase the no. of iterations (curve goes down)
  - Very useful in debugging and seeing the behavior of online learning algorithms
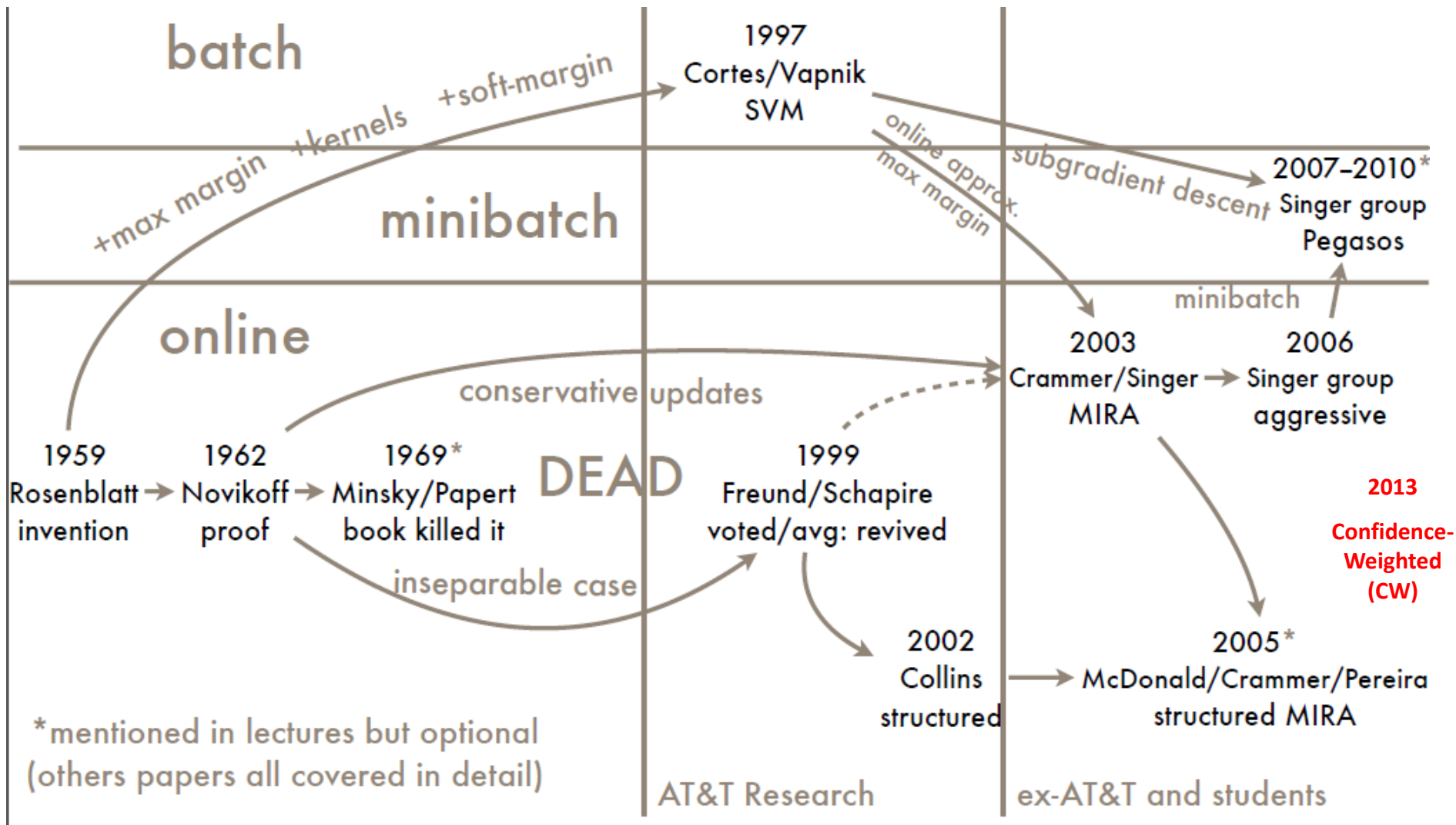
- **Hyper-parameter Optimization**
  - Split the training data: sub-train + validation data
  - Tune hyper-parameters (e.g., no. of iterations) on the validation data
  - The learner should not look at the test data!

# Two kinds of learning curves

- General learning curve
  - Applicable for any learning algorithm
  - X-axis (number of training examples)
  - Y-axis (accuracy on unseen data)

- Online learning curve
  - Specific to online learning algorithms
  - For a fixed training set
  - X-axis (number of iterations or passes over the data)
  - Y-axis (number of mistakes made)

# History of Perceptron*



* slide from Liang Huang

41

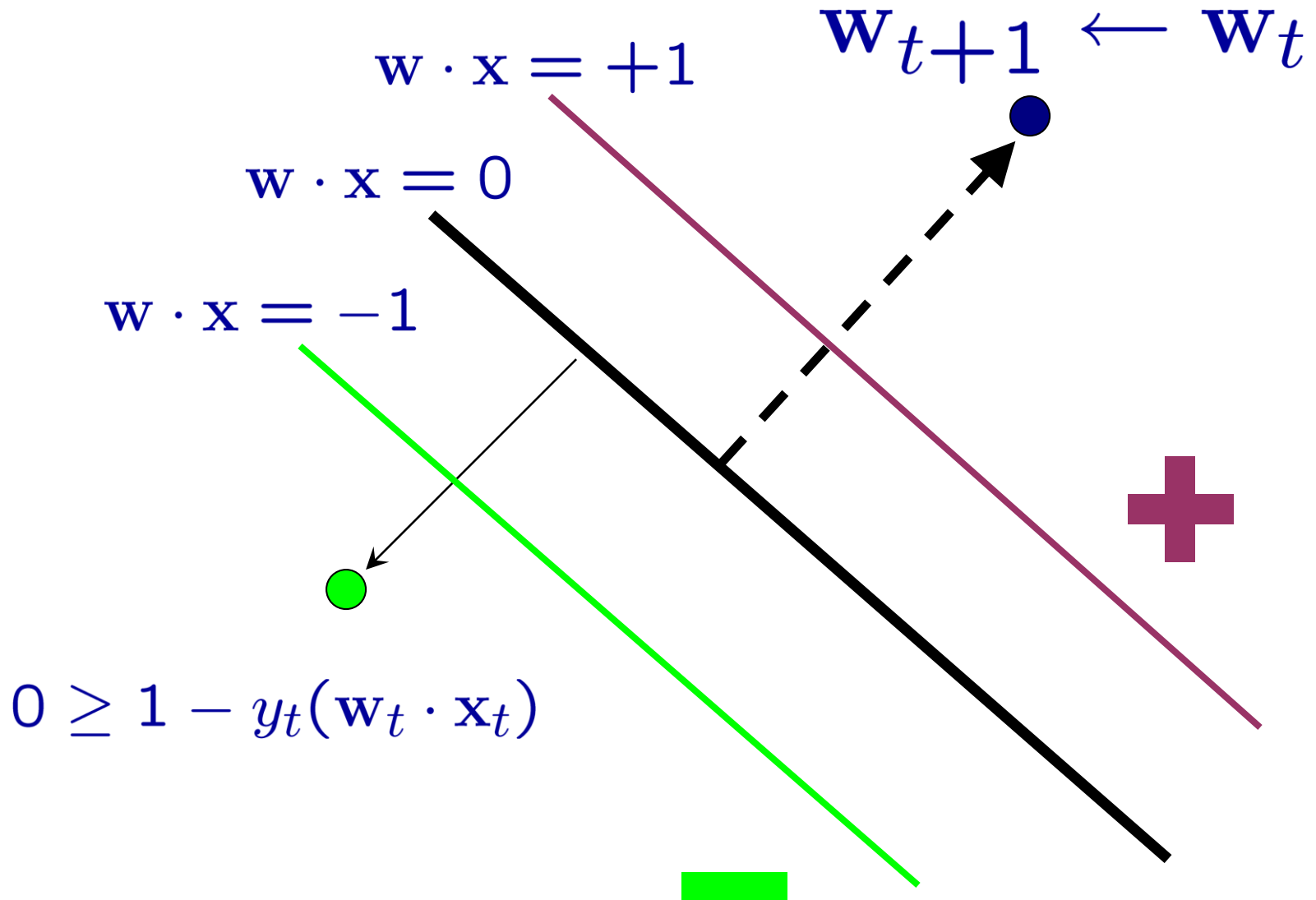# Passive-Aggressive (PA) Algorithm

# PA Algorithm: Motivation

- **Perceptron:** No guaranties of margin *after* the update

- **PA:** Enforce a minimal non-zero margin after the update

- In particular :
  - If the margin is large enough (1), then do nothing
  - If the margin is less then unit, update such that the margin *after* the update is *enforced* to be unit

# Margin

- Y.(w_{t+1}.x) > 0 then the prediction is correct

  - +1 ( > 0) > 0 => Sign (w_{t+1}.x) = +1

  - -1 (< 0) > 0 => Sign (w_{t+1}.x) = -1

- If I'm making a mistake, then this margin is < 0

- If margin > 0 implies corrective behavior

  - Margin > 1 also implies corrective behavior

# Input Space

$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$

$\mathbf{w} \cdot \mathbf{x} = +1$

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

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

$0 \geq 1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)$

# Input Space vs. Version Space

- **Input Space :**
  - Points are input data $y_t \mathbf{x}_t$
  - One constraint is induced by weight vector $\mathbf{w}$
  - Primal space
  - Half space = all input examples that are classified correctly by a given predictor (weight vector)
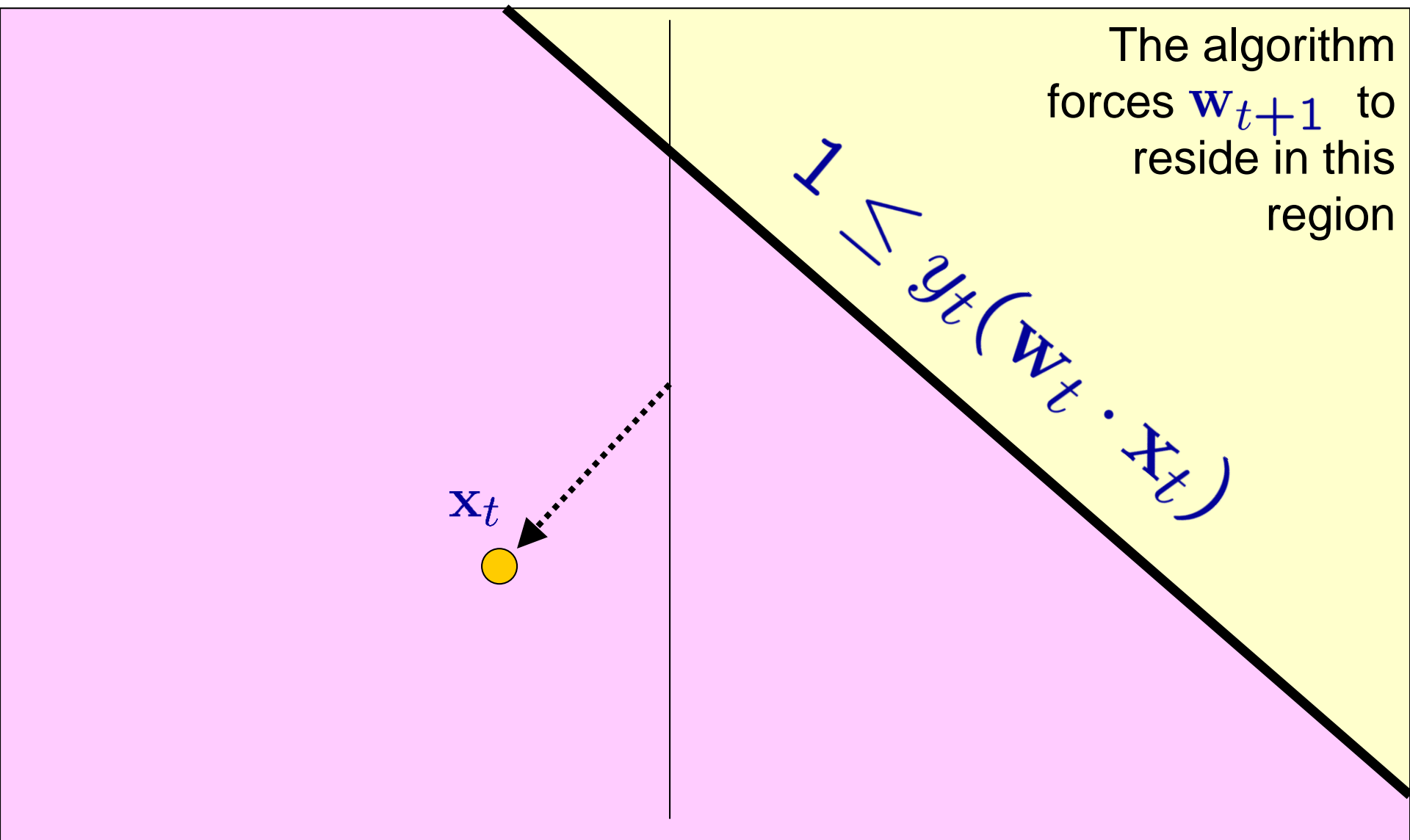
$$\{ y\mathbf{x} \ : \ \mathbf{w} \cdot (y\mathbf{x}) \geq 0 \}$$

- **Version Space :**
  - Points are weight vectors $\mathbf{w}$
  - One constraints is induced by input data $y_t \mathbf{x}_t$
  - Dual space
  - Half space = all predictors (weight vectors) that classify correctly a given input example

$$\{ \mathbf{w} \ : \ \mathbf{w} \cdot (y\mathbf{x}) \geq 0 \}$$

# Weight vector (Version) Space



The algorithm forces $\mathbf{w}_{t+1}$ to reside in this region

$$1 \leq y_t(\mathbf{w}_t \cdot \mathbf{x}_t)$$

$\mathbf{x}_t$

# Passive Step

$$1 \leq y_t(\mathbf{w}_t \cdot \mathbf{x}_t)$$

$\mathbf{w}_{t+1}$

Nothing to do.

$\mathbf{w}_t$ already resides on the desired side.

# Aggressive Step



$1 \leq y_t(\mathbf{w}_t \cdot \mathbf{x}_t)$

$\mathbf{w}_t$

$\mathbf{w}_{t+1}$

The algorithm projects $\mathbf{w}_t$ on the desired half-space

# Aggressive Update Step

- Set $\mathbf{w}_{t+1}$ to be the solution of the following optimization problem :

$$\mathbf{w}_{t+1} = \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 \quad \text{\textbf{Conservative}}$$
$$\text{s.t.} \quad y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq 1 \quad \text{\textbf{Corrective}}$$

- The Lagrangian :

$$\mathcal{L}(\mathbf{w}, \tau) = \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 + \tau(1 - y_t(\mathbf{w} \cdot \mathbf{x}_t))$$

- Solve for the dual : $\quad \max_{\tau \geq 0} \min_{\mathbf{w}} \quad \mathcal{L}(\mathbf{w}, \tau)$

# Norm of a vector

- L_1 norm

- L_2 norm

- L_infinity norm

- W = [w_1, w_2,...,w_d]

- ||W|| = Sqrt (w_1^2 + w_2^2 +...+w_d^2)

# Aggressive Update Step

- Optimize for $\mathbf{w}$ :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \mathbf{w}_t - \tau y_t \mathbf{x}_t$$

- Set the derivative to zero $\quad \mathbf{w} = \mathbf{w}_t + \tau y_t \mathbf{x}_t$

- Substitute back into the Lagrangian :

$$\mathcal{L}(\tau) = -\frac{1}{2}\|\mathbf{x}_t\|^2 \tau^2 + \tau(1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t))$$

- Dual optimization problem

$$\max_{\tau \geq 0} \quad -\frac{1}{2}\|\mathbf{x}_t\|^2 \tau^2 + \tau(1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t))$$

# Aggressive Update Step

- Dual Problem :

$$\max_{\tau \geq 0} \quad -\frac{1}{2}\|\mathbf{x}_t\|^2 \tau^2 + \tau(1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t))$$

- Solve it :

$$\tau = \max\left\{0, \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}\right\}$$

# Alternative Derivation

- Additional Constraint (linear update) :

$$\mathbf{w}_{t+1} \;=\; \mathbf{w}_t + \tau y_t \mathbf{x}_t$$

- Force the constraint to hold as equality

$$1 = y_t((\mathbf{w}_t + \tau y_t \mathbf{x}_t) \cdot \mathbf{x}_t) \;=\; y_t(\mathbf{w}_t \cdot \mathbf{x}_t) + \tau \|\mathbf{x}_t\|^2$$
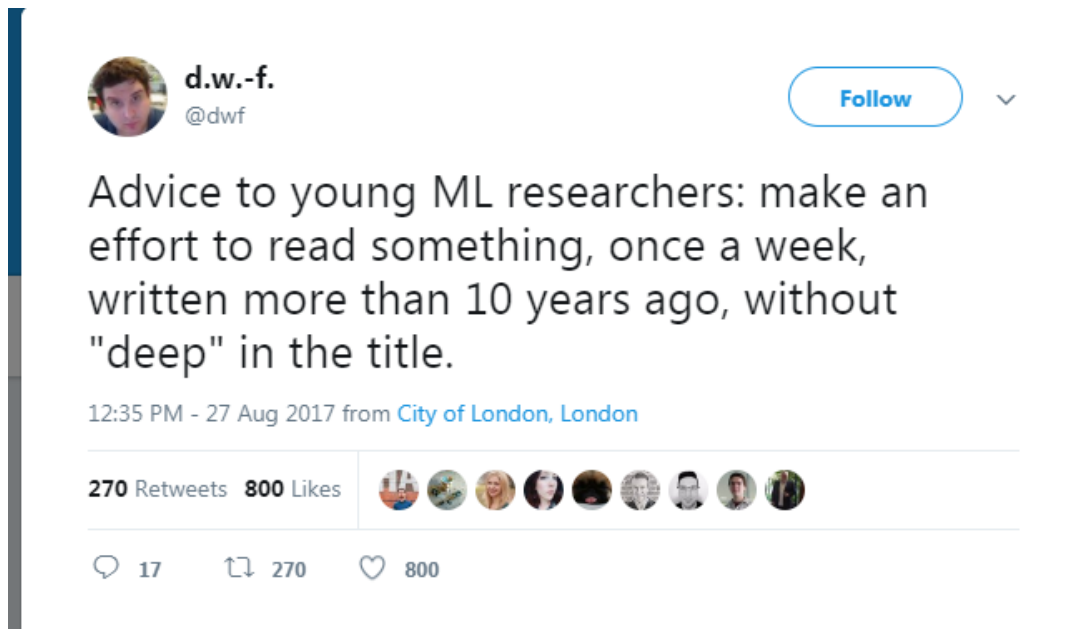
- Solve :

$$\tau = \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}$$
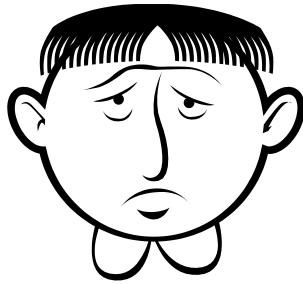
# Recap of last Lecture

- Perceptron algorithm
  - ▲ Convergence result for linearly separable data
  - ▲ Voted and Averaged perceptron
  - ▲ Practical tricks: shuffling, variable learning rate, …

- History of Perceptron
  - ▲

# Recap of Last Lecture

- **Passive-Aggressive Algorithm**
  - Smallest change to weights to ensure a margin of 1 on the new example
  - Derivation of weight update
  - Perceptron vs. Passive-Aggressive update

# Passive-Aggressive Update





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

$$y_t(\mathbf{w}_t \cdot \mathbf{x}_t) \geq 1 \qquad\qquad y_t(\mathbf{w}_t \cdot \mathbf{x}_t) < 1$$

$$\tau = 0 \qquad\qquad \tau = \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}$$

# Perceptron vs. PA Update

- Common Update :

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

- Perceptron

$$\tau = \begin{cases} 1 & y_t(\mathbf{w}_t \cdot \mathbf{x}_t) > 0 \\ 0 & y_t(\mathbf{w}_t \cdot \mathbf{x}_t) \leq 0 \end{cases}$$

- Passive-Aggressive

$$\tau = \max\left\{0, \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}\right\}$$
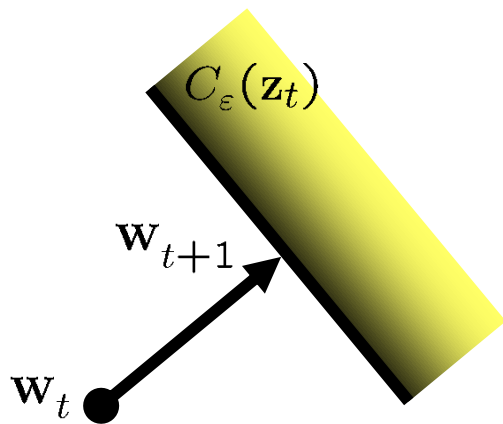
# Two kinds of updates

- Additive update
  - $w_{t+1} = w_t + \delta$ (adding some gradient)

- Multiplicative update
  - $W_{t+1} = w_t \times \delta$
  - Example: Exponentiated Gradient Algorithm

- Yoram Singer => Hebrew Univ => Google (2007)
  - Spam classifiers in Gmail => they are passive-aggressive algorithm
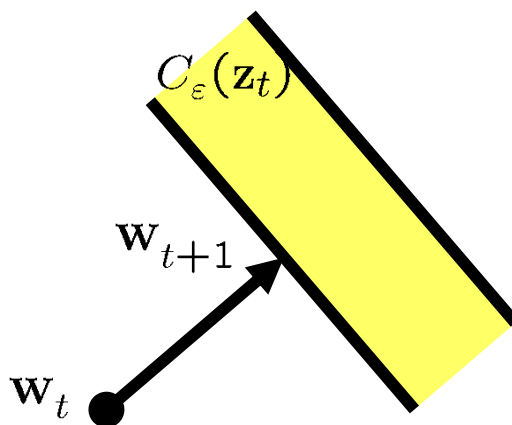
# The Passive-Aggressive Algorithm

- Each example defines a set of consistent hypotheses: $C_\varepsilon(\mathbf{z}_t) = \{\mathbf{w} \mid \delta(\mathbf{w}; \mathbf{z}_t) \le \varepsilon\}$

- The new vector $\mathbf{w}_{t+1}$ is set to be the projection of $\mathbf{w}_t$ onto $C_\varepsilon(\mathbf{z}_t)$

$$\mathbf{w}_{t+1} = \arg\min_{\mathbf{w}} \|\mathbf{w} - \mathbf{w}_t\| \ \ \text{s.t.} \ \ \mathbf{w} \in C_\varepsilon(\mathbf{z}_t)$$

Classification          Regression



$C_\varepsilon(\mathbf{z}_t)$          $C_\varepsilon(\mathbf{z}_t)$

$\mathbf{w}_{t+1}$          $\mathbf{w}_{t+1}$

$\mathbf{w}_t$          $\mathbf{w}_t$

# Unrealizable Case

There is no weight vector that satisfy all the constraints

# Unrealizable Case

$$\mathbf{w}_{t+1} = \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi^2$$

$$\text{s.t.} \quad y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq 1 - \xi \qquad \xi \geq 0$$

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

$\xi$

$\xi^2$

$$\min\left\{C, \max\left\{0, \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}\right\}\right\}$$

$$\max\left\{0, \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2 + C}\right\}$$

62

# Reduction Hammer

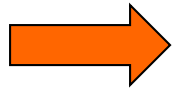- Smaller problem/nail (binary classification) => Simpler solution/hammer (Perceptron or PA)
  - You have only two classes

- Harder problem/nail
  - You have more than two classes
  - How can I develop an appropriate hammer for this nail?

- Create some simpler problems that you know how to solve => Solve them => Aggregate the solutions to come up with a solution for the harder problem

# Reduction for Multi-Class Classification
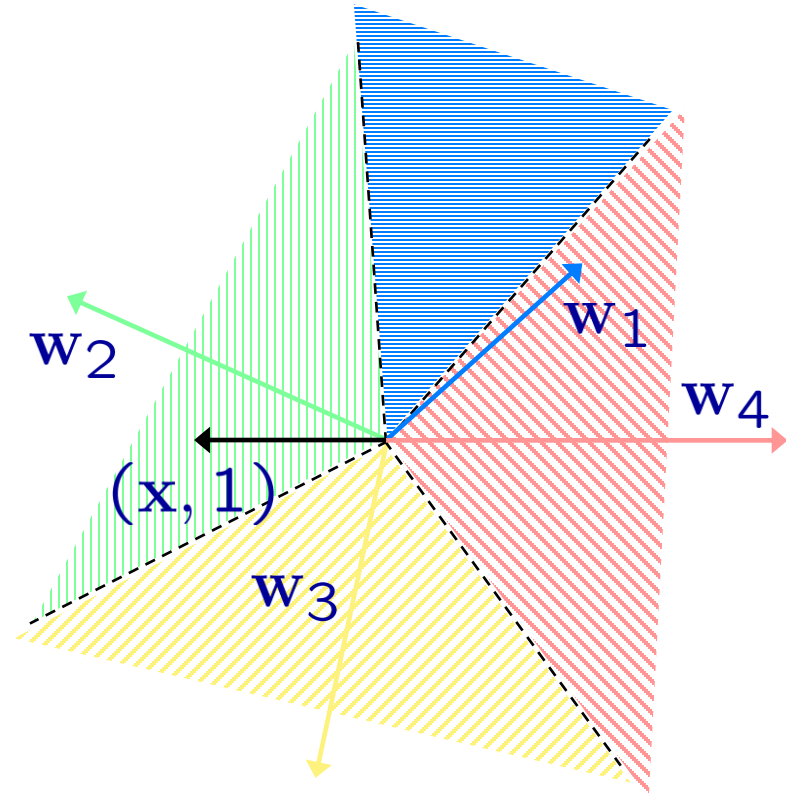
- Classes = {red, blue, green}

- Bunch of training examples {x,y} pairs where y is one of the three classes

- One against others {red vs. other; blue vs. others; green vs. others}

- Increasing the dimension of margin

- Tree of classifiers to filter candidate class labels

- One vs. One: Red vs. Blue; Blue vs. Green; Green vs. Red => Given a new example => Majority vote

# Multi-Class: Representation-I

- k Prototypes $\mathbf{w}_1, \mathbf{w}_2 \ldots \mathbf{w}_k$

- New instance $\mathbf{x}$

- Compute $\mathrm{Score}(r) = \mathbf{w}_r \cdot \mathbf{x}$

| Class r | $\mathbf{w}_r \cdot \mathbf{x}$ |
|---------|--------------------------------|
| 1 | -1.08 |
| 2 | 1.66 |
| 3 | 0.37 |
| 4 | -2.09 |



- Prediction:

  The class achieving the highest Score

# Multi-Class Representation-II

- Weight-vector per class (Representation I)
  - ▲ Intuitive

- Single weight-vector (Representation II)
  - ▲ Generalizes representation I

F(x,4) =

| 0 | 0 | 0 | x | 0 |
|---|---|---|---|---|

- Predict label with highest score (Inference)

$$\arg \max_{\mathbf{z}} F(\mathbf{x}, \mathbf{z}) \cdot \mathbf{w}$$

# Margin for Multi-Class

- Binary :

$$\mathbf{w}_2 \cdot \mathbf{x} - \mathbf{w}_1 \cdot \mathbf{x} \geq 1$$

$$\mathbf{w} \cdot F(\mathbf{x}, 2) - \mathbf{w} \cdot F(\mathbf{x}, 1) \geq 1$$

- Multi Class :

$$\mathbf{w}_y \cdot \mathbf{x} - \mathbf{w}_z \cdot \mathbf{x} \geq 1 \quad \forall z \neq y$$

$$\mathbf{w} \cdot F(\mathbf{x}, y) - \mathbf{w} \cdot F(\mathbf{x}, z) \geq 1$$
$$\forall z \neq y$$

# Margin for Multi-Class

- Multi Class :

$$\mathbf{w} \cdot F(\mathbf{x}, y) - \mathbf{w} \cdot F(\mathbf{x}, z) \geq 1$$
$$\forall z \neq y$$

Because the loss function is not constant !
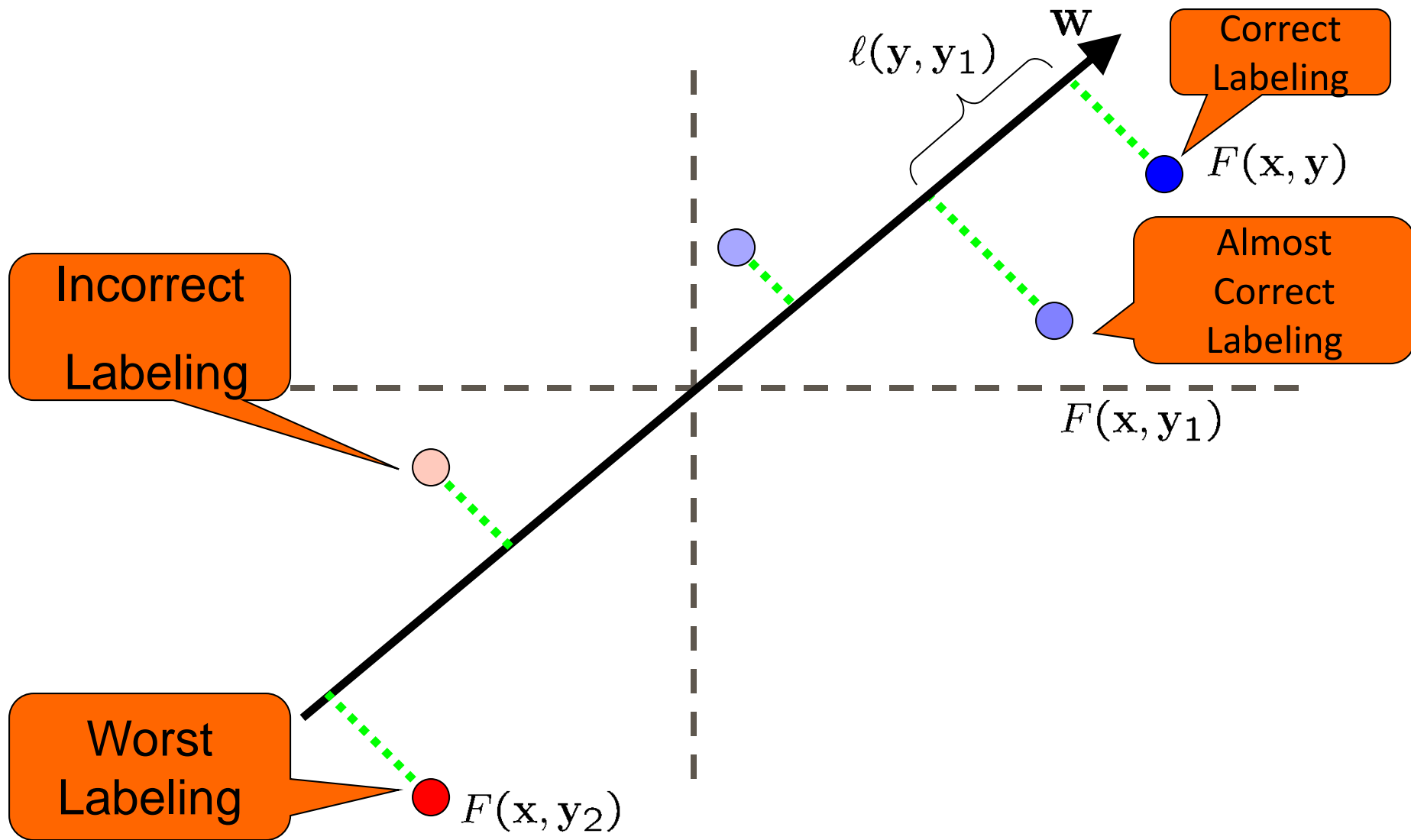
# Margin for Multi-Class

- Multi Class :

$$\mathbf{w} \cdot F(\mathbf{x}, y) - \mathbf{w} \cdot F(\mathbf{x}, z) \geq \ell(z, y)$$
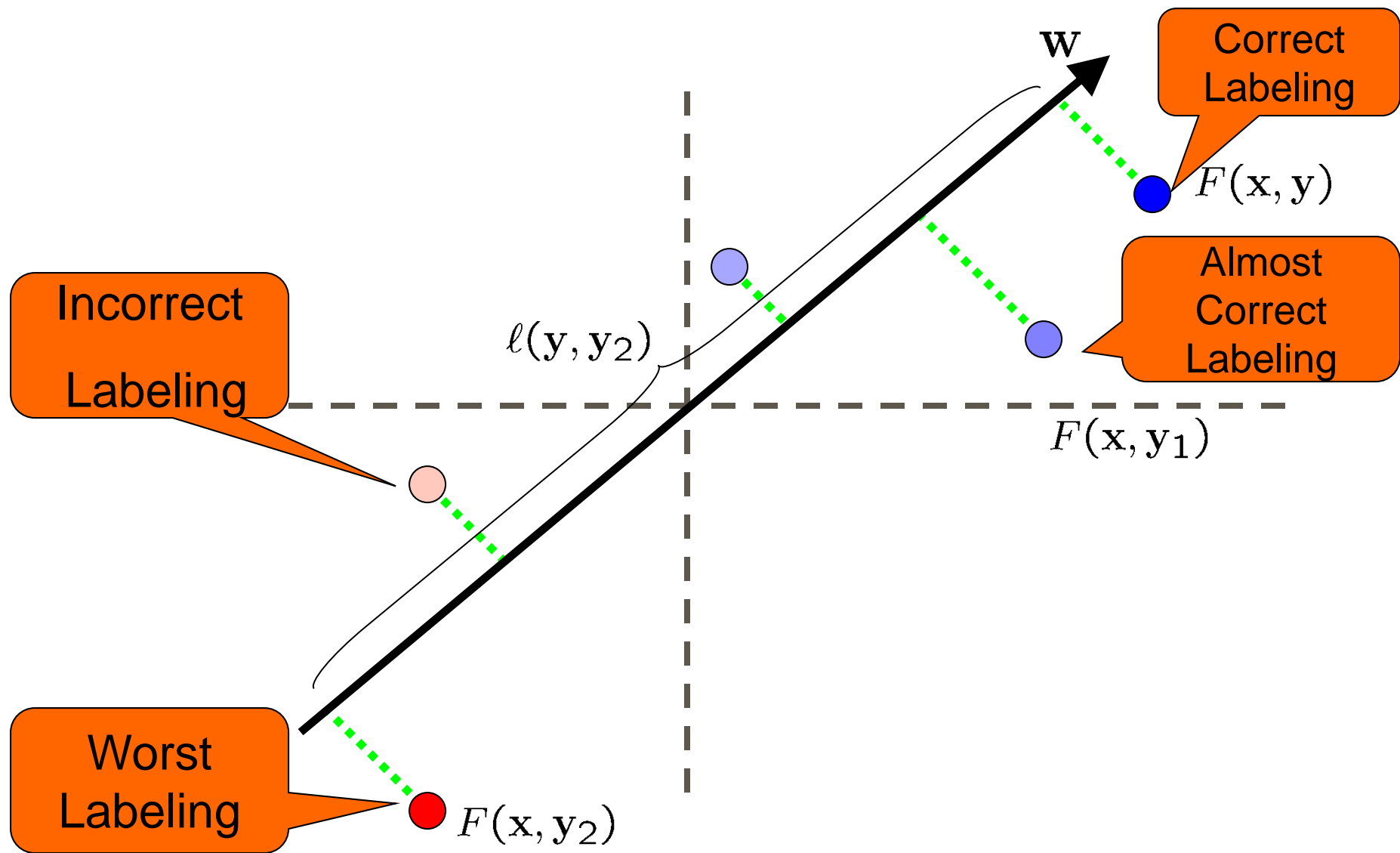$$\forall z \neq y$$

So, use it !

# Margin Scaled by Loss: Illustration

# Margin Scaled by Loss: Illustration

# PA Multi-Class Update

- Project the current weight vector such that the instance ranking is consistent with loss function

- Set $\mathbf{w}_{t+1}$ to be the solution of the following optimization problem :

$$\begin{aligned}
\mathbf{w}_{t+1} \;=\; & \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 \\
\text{s.t.} \quad & \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, z) \geq \ell(z, y_t) \\
& \forall z \neq y_t
\end{aligned}$$

# PA Multi-Class Update

- **Problem**
  - intersection of constraints may be empty

$$\left\{ \mathbf{w} \; : \; \begin{array}{c} \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, z) \geq \ell(z, y_t) \\ \forall z \neq y_t \end{array} \right\} = \emptyset$$

- **Solutions**
  - Does not occur in practice
  - Add a slack variable
  - Remove constraints

# Add a Slack Variable

- Add a slack variable :

$$\mathbf{w}_{t+1} = \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi$$

$$\text{s.t.} \quad \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, z) \geq \ell(z, y_t) - \xi$$

$$\forall z \neq y_t$$

$$\xi \geq 0$$

- Rewrite the optimization :

$$\frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 + C \max \left\{ \begin{array}{c} \mathbf{w} \cdot F(\mathbf{x}_t, z) - \mathbf{w} \cdot F(\mathbf{x}_t, y_t) + \ell(z, y_t) \\ 0 \end{array} \right\}$$

Generalized Hinge loss!!

# PA Multi-Class Update

- Remove constraints :

$$\mathbf{w}_{t+1} = \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2$$

$$\text{s.t.} \quad \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, z) \geq \ell(z, y_t)$$

$$\forall z \neq y_t$$

$$\mathbf{w}_{t+1} = \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2$$

$$\text{s.t.} \quad \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, \widehat{y}_t) \geq \ell(\widehat{y}_t, y_t)$$

- How to choose the single competing labeling?

  - The labeling that attains the highest score!

$$\widehat{y}_t = \arg\max_z \ \mathbf{w}_t \cdot F(\mathbf{x}_t, z)$$

  - … which is the predicted label according to the current model

# PA Multi-Class Online Algorithm

- Initialize $\mathbf{w}_1$

- For $t = 1 \ldots T \ldots$
  - Receive an input instance $\mathbf{x}_t$
  - Outputs a prediction $\widehat{y}_t = \arg\max_z \; \mathbf{w}_t \cdot F(\mathbf{x}_t, z)$
  - Receives a feedback label $y_t$
  - Computes loss $\ell(\widehat{y}_t, y_t)$
  - Update the prediction rule

$$
\begin{aligned}
\mathbf{w}_{t+1} \;=\; & \min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_t\|^2 \\
& \text{s.t.} \quad \mathbf{w} \cdot F(\mathbf{x}_t, y_t) - \mathbf{w} \cdot F(\mathbf{x}_t, \widehat{y}_t) \geq \ell(\widehat{y}_t, y_t)
\end{aligned}
$$

# Multi-Class Online Algorithm

- Initialize   $\mathbf{w}_1$

- For   $t = 1 \ldots T \ldots$

    - Receive an input instance   $\mathbf{x}_t$
    - Outputs a prediction   $\hat{y}_t = \arg\max_z \ \mathbf{w}_t \cdot F(\mathbf{x}_t, z)$
    - Receives a feedback label   $y_t$
    - Computes loss   $\ell(\hat{y}_t, y_t)$
    - Update the prediction rule

$$w_{t+1} = w_t + \tau \cdot (F(x_t, y_t) - F(x_t, \hat{y}_t))$$

# Recap of Last Lecture

- **Multi-Class Representations**
  - Multi-prototype (one weight vector $\boldsymbol{w_i}$ for each class $\boldsymbol{i}$)
  - Single-prototype (one weight vector $\boldsymbol{w}$ – concatenation of all the $\boldsymbol{k}$ weight vector) via extended feature space $\boldsymbol{F(x, y)}$

- **Multi-Class Passive-Aggressive Algorithm**
  - Margin for multi-class classification
  - Margin scaled by the loss function
  - Mathematical optimization problem with $\boldsymbol{k-1}$ constraints

# Multi-Class Online Algorithm

- Initialize $\mathbf{w}_1$
- For $t = 1 \ldots T \ldots$
  - Receive an input instance $\mathbf{x}_t$
  - Outputs a prediction $\hat{y}_t = \arg\max_z \; \mathbf{w}_t \cdot F(\mathbf{x}_t, z)$
  - Receives a feedback label $y_t$
  - Computes loss $\ell(\hat{y}_t, y_t)$
  - Update the prediction rule

$$w_{t+1} = w_t + \tau \cdot (F(x_t, y_t) - F(x_t, \hat{y}_t))$$

$$\tau = \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}$$

$$\tau = \frac{1 - (w_t \cdot F(x_t, y_t) - w_t \cdot F(x_t, \hat{y}_t))}{\|F(x_t, y_t) - F(x_t, \hat{y}_t)\|^2}$$

**Binary**                    **Multi-class**        80

# Multi-Class Classification Implementation

- "d" features for each input example

- W1 (d weights) for class 1

- W2 (d weights for class 2

- ..

- Score(class 1) = W1.X

- W_{correct} = W_{correct} + \eta.X

- W_{wrong} = W_{wrong} - \eta.X

# Multi-Class Update

- Score(correct-label) = 3.4
  - W_{correct-label}.X

- Score(wrong-label) = 5.6
  - W_{wrong-label}.X

- Weight update
  - W_{correct-label} = W_{correct-label} + \tau. X
  - W_{wrong-label} = W_{wrong-label} - \tau.X

- Updated score(correct-label)
  - (W_{correct-label} + \tau. X).X = \tau.X^2 (added)

- Updated score(wrong-label)

# Binary vs. Multi-Class Interpretation

- Correct-label $y_t$

- Wrong-label $\hat{y_t}$

- Feature vectors
  - $F(x, y_t)$
  - $F(x, \hat{y_t})$

- $F(x, y_t) - F(x, \hat{y_t})$ == Input feature vector

- Output label = +1

- $F(x, \hat{y_t}) - F(x, y_t)$ == Input feature vector

- Output label = -1

# Confusion Matrix

- Measures which classes are easy or hard to separate -- $k$ classes implies $k \times k$ matrix

$$k$$

$$k$$



- $C_{ij} = \dfrac{\# \ examples \ with \ class \ label \ j \ that \ are \ classified \ as \ label \ i}{\# \ examples \ with \ label \ j}$

- High diagonal entry implies class is easy to classify

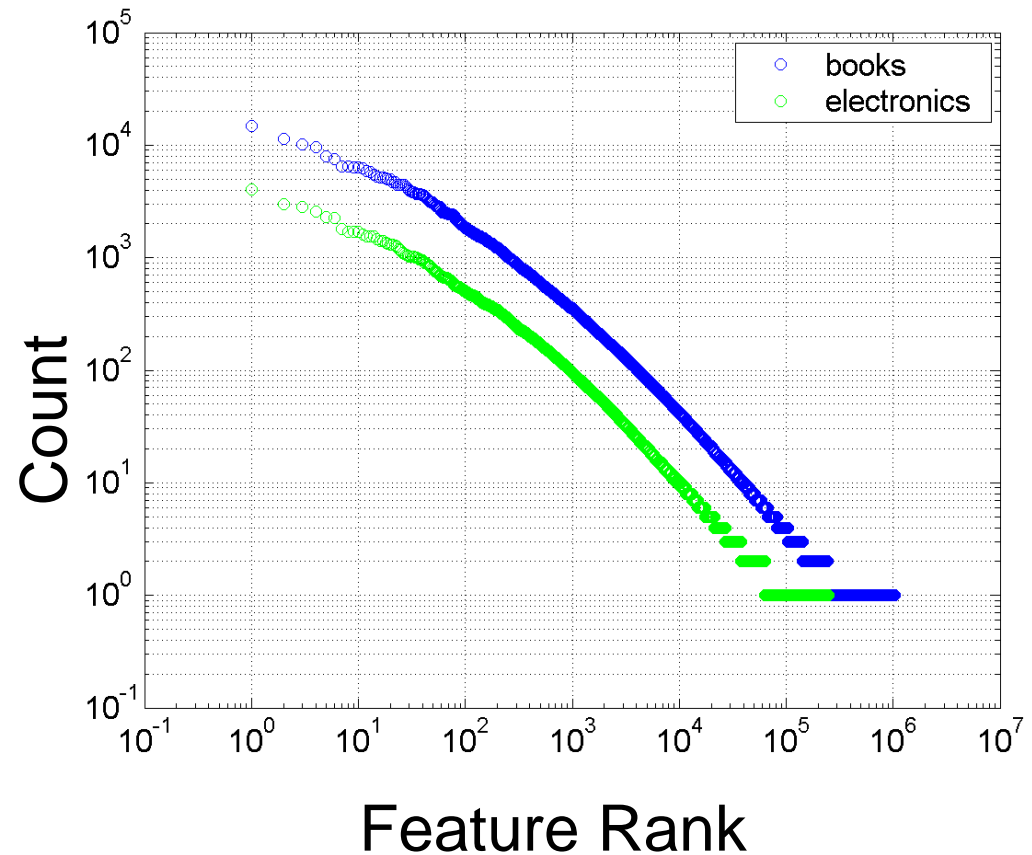- High off-diagonal entry implies classes are easily confused

# Learning Curves: Online vs. General

- **Online learning curve**
  - training iterations vs. Mistakes
  - only applicable for online learning algorithms

- **General learning curve**
  - amount of training data vs. accuracy
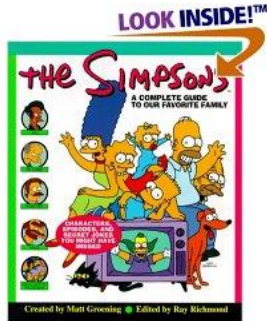  - applicable to any learning algorithm

# Confidence-Weighted (CW) Algorithm

# Motivation: NLP problems

- Big datasets, large number of features

- Many features are only weakly correlated with target label

- Heavy-tailed feature distribution

- Linear classifiers: features are associated with word counts

# Motivation: Sentiment Classification



- **Who needs this Simpsons book? You DOOOOOOOO**

  This is one of the most extraordinary volumes I've ever encountered encapsulating a television series … . Exhaustive, informative, and ridiculously entertaining, it is the best accompaniment to the best television show … . Even if you only "enjoy" the Simpsons (as opposed to being a raving fanatic, which most people who watch the show are, after all … Very highly recommended!

# Motivation: Sentiment Classification

- Many positive reviews with the word best

$$W_{best}$$

- Later negative review

  - "*boring book – best if you want to sleep in seconds*"

- Linear update will reduce both

$$W_{best} \quad W_{boring}$$

- But best appeared more than boring

- How to adjust weights at different rates?  $W_{boring}$  $W_{best}$

# Span based Update Rules

- The weight vector is a linear combination of examples

$$w_f \leftarrow w_f + \eta \, y \, x_f \quad \forall f$$

Weight of feature f

Learning rate

Target label, -1 or 1

Value of feature f of instance
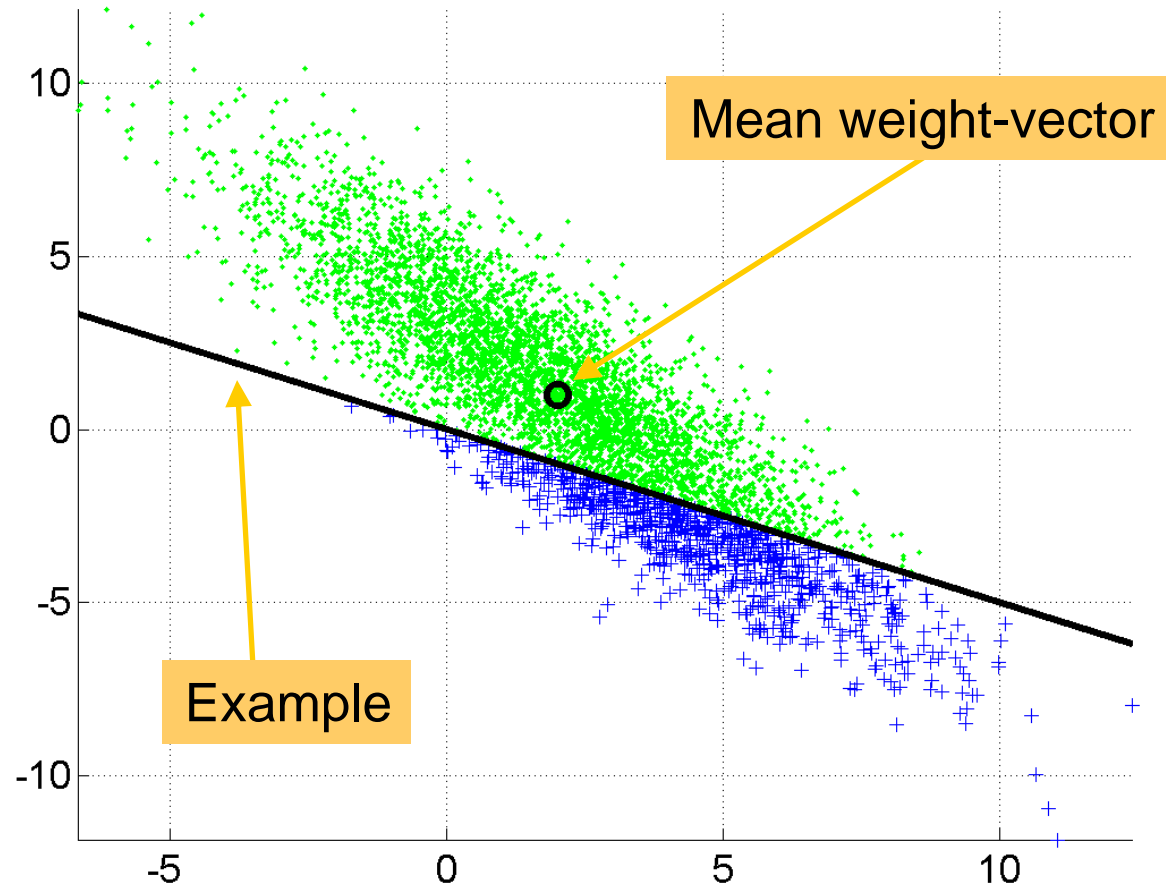
- Two rate schedules (among others):
  - Perceptron algorithm, conservative:
  
  $$\eta = 1$$
  
  - Passive-aggressive
  
  $$\eta = \max \left\{ 0, \frac{1 - y(\mathbf{w} \cdot \mathbf{x})}{\|\mathbf{x}\|^2} \right\}$$

# Distributions in Version Space



Mean weight-vector

Example

# Margin as a Random Variable

- Signed margin

$$M = y(\boldsymbol{w} \cdot \boldsymbol{x})$$

is a Gaussian-distributed variable

$$M \sim \mathcal{N}\left(y(\boldsymbol{x} \cdot \boldsymbol{\mu}) \ , \ \boldsymbol{x}^\top \Sigma \boldsymbol{x}\right)$$

- Thus:

$$\Pr\left[y\left(\boldsymbol{w} \cdot \boldsymbol{x}\right) \geq 0\right] = \Phi\left(\frac{y(\boldsymbol{x} \cdot \boldsymbol{\mu})}{\sqrt{\boldsymbol{x}^\top \Sigma \boldsymbol{x}}}\right)$$

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{z} e^{-t^2} dt$$

# Version (weight vector) Space



Place most of the probability mass in this region
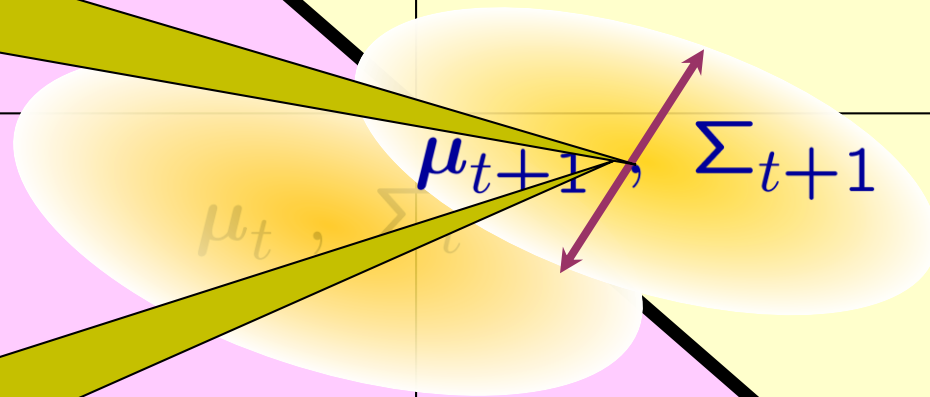
$\mathbf{w}$

$\mathbf{x}_t$

$0 \leq y_t(\mathbf{w} \cdot \mathbf{x}_t)$

# Aggressive Step

Mean moved past the mistake line

(large margin)

The covariance is shirked in the direction of the new example

$\boldsymbol{\mu}_{t+1}, \; \boldsymbol{\Sigma}_{t+1}$

$\mu_t \; , \; \Sigma_t$

# PA like Update

- PA:

$$\min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w} - \mathbf{w}_i\|^2$$
$$\text{s.t.} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i) \geq 1$$

- New Update :

$$\min_{\boldsymbol{\mu},\Sigma} \quad D_{\mathsf{KL}}\left(\mathcal{N}\left(\boldsymbol{\mu},\Sigma\right) \| \mathcal{N}\left(\boldsymbol{\mu}_i,\Sigma_i\right)\right)$$
$$\text{s.t.} \quad \Pr\left[y_i\left(\boldsymbol{w}\cdot\boldsymbol{x}_i\right) \geq 0\right] \geq \eta$$

Confidence Parameter
$$\eta \in (0.5, 1)$$

# Summary of Online Learning

- **Online learning**
  - Iterative game between teacher and learner

- **Design principles of online learning**
  - Trade-off amount of change (conservative) and reduction in loss (corrective)

- **Online learning algorithms**
  - Perceptron (fixed learning rate for all examples)
  - Passive-Aggressive (fixed learning rate for each example)
  - Confidence-weighted classifier (fixed learning for each feature and each example)

# Questions?