	Introduction
	 Sometimes probabilistic information unavailable or mathematically intractable
CSCE 970 Lecture 3: Linear Classifiers	• Many alternatives to Bayesian classification, but optimality guarantee may be compromised!
	 <u>Linear classifiers</u> use a <u>decision hyperplane</u> to perform classification
Stephen D. Scott	• Simple and efficient to train and use
	 Optimality requires <u>linear separability</u> of classes
January 22, 2001	$\bigcirc = \text{Class A} + = \text{Class B} \\ & & \\ & \\ & \\ & \\ & \\ & \\ & \\ & \\ & $
1	2
	The Perceptron Algorithm
Linear Discriminant Functions	• Assume linear separability, i.e. $\exists w^*$ s.t.
• Let $\mathbf{w} = [w_1, \dots, w_\ell]^T$ be a <u>weight vector</u> and w_0 (a.k.a. θ) be a <u>threshold</u>	$ \begin{split} \mathbf{w}^{*T} \cdot \mathbf{x} &> 0 \forall \mathbf{x} \in \omega_1 \\ \mathbf{w}^{*T} \cdot \mathbf{x} &\leq 0 \forall \mathbf{x} \in \omega_2 \\ (w_0^* \text{ is included in } \mathbf{w}^*) \end{split} $
• Decision surface is a hyperplane: $\mathbf{w}^T \cdot \mathbf{x} + w_0 = 0$	 So ∃ deterministic function clasifying vectors (contrary to Ch. 2 assumptions)
• E.g. predict ω_2 if $\sum_{i=1}^\ell w_i x_i > w_0,$ otherwise predict ω_1	$x_{1} \xrightarrow{w_{1}} \hat{y}(t) = 0 \text{ otherwise}$ (ω_{1}) (ω_{1}) $\hat{y}(t) = 1 \text{ if sum} > w_{0}$ $\hat{y}(t) = 0 \text{ otherwise}$
• Focus of this lecture: How to find w_i 's	(ω_2) May also use +1 and -1
 Perceptron algorithm 	• Given actual label $y(t)$ for <u>trial</u> t , update weights:
– Winnow	$\mathbf{w}(t+1) = \mathbf{w}(t) + \rho(\mathbf{y}(t) - \hat{\mathbf{y}}(t))\mathbf{x}(t)$
 Least squares methods (if classes not lin- early separable) 	• $\rho > 0$ is learning rate • $(y(t) - \hat{y}(t))$ moves weights toward correct prediction for x



	Winnow/Exponentiated Gradient Intuition
Winnow/Exponentiated Gradient $ \begin{array}{c} & \begin{array}{c} & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & $	• Measure distance in cost function with <u>unnormalized relative entropy</u> : $U(\mathbf{w}) = \sum_{i=1}^{\ell} \left(w_i(t) - w_i(t+1) + w_i(t+1) \ln \frac{w_i(t+1)}{w_i(t)} \right)$ $+ \widehat{\eta} (y - \mathbf{w}(t+1) + \mathbf{w}(t+1) \ln \frac{w_i(t+1)}{w_i(t)} + \widehat{\eta} (y - \mathbf{w}(t+1) + \mathbf{w}(t+1))^2$ • Take gradient w.r.t. $\mathbf{w}(t+1)$ and set to 0: $0 = \ln \frac{w_i(t+1)}{w_i(t)} - 2\eta \left(y(t) - \sum_{i=1}^{\ell} w_i(t+1)x_i(t) \right) x_i(t)$ • Approximate with $0 = \ln \frac{w_i(t+1)}{w_i(t+1)} - 2\eta \left(y(t) - \sum_{i=1}^{\ell} w_i(t)x_i(t) \right) x_i(t)$
$w_i(t)$ if $\hat{y}(t) = y(t)$	$w_{i}(t) = 2\eta \left(y(t) - \sum_{i=1}^{n} w_{i}(t) x_{i}(t) \right) x_{i}(t),$ which yields $w_{i}(t+1) = w_{i}(t) \exp\left(-2\eta \left(\hat{y}(t) - y(t)\right) x_{i}(t)\right)$ 10
Winnow/Exponentiated Gradient Negative Weights	Winnow/Exponentiated Gradient Miscellany
 Winnow/Exponentiated Gradient Negative Weights Winnow and EG update wts by mult by a pos const, impossible to change sign Woight vectors restricted to one guadrant 	 Winnow/Exponentiated Gradient Miscellany Winnow and EG are <u>muliplicative weight update</u> schemes versus <u>additive weight update</u> schemes, e.g. Perceptron
 Winnow/Exponentiated Gradient Negative Weights Winnow and EG update wts by mult by a pos const, impossible to change sign Weight vectors restricted to one quadrant Solution: Maintain wt vectors w⁺(t) and w⁻(t) 	 Winnow/Exponentiated Gradient Miscellany Winnow and EG are <u>muliplicative weight update</u> schemes versus <u>additive weight update</u> schemes, e.g. Perceptron Winnow and EG work well when most attributes (features) are <u>irrelevant</u>, i.e. optimal weight vector w* is sparse (many 0 entries)
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update:	 Winnow/Exponentiated Gradient Miscellany Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries) E.g. x_i ∈ {0, 1}, x's are labelled by a monotone k-disjunction over l attributes, k ≪ l
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update: $r_i^+(t) = \exp(-2\eta (\hat{y}(t) - y(t)) x_i(t)U)$	 Winnow/Exponentiated Gradient Miscellany Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries) E.g. x_i ∈ {0, 1}, x's are labelled by a monotone k-disjunction over l attributes, k ≪ l Remaining l – k are irrelevant
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos- const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update: $r_i^+(t) = \exp(-2\eta (\hat{y}(t) - y(t)) x_i(t) U)$ $r_i^-(t) = 1/r_i^+(t)$	 Winnow/Exponentiated Gradient Miscellany Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries) E.g. x_i ∈ {0, 1}, x's are labelled by a monotone k-disjunction over ℓ attributes, k ≪ ℓ Remaining ℓ - k are irrelevant E.g. x₅ ∨ x₉ ∨ x₁₂, ℓ = 150, k = 3
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update: $r_i^+(t) = \exp(-2\eta (\hat{y}(t) - y(t)) x_i(t)U)$ $r_i^-(t) = 1/r_i^+(t)$ $w_i^+(t+1) = U \cdot \frac{w_i^+(t)r_i^+(t)}{\sum_{j=1}^{\ell} (w_i^+(t)r_i^+(t) + w_i^-(t)r_i^-(t))}$	Winnow/Exponentiated Gradient Miscellany• Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron• Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries)• E.g. $x_i \in \{0, 1\}$, x's are labelled by a monotone k -disjunction over ℓ attributes, $k \ll \ell$ - Remaining $\ell - k$ are irrelevant- E.g. $x_5 \lor x_9 \lor x_{12}$, $\ell = 150$, $k = 3$ - For disjunctions, number of on-line prediction mistakes is $O(k \log \ell)$ for Winnow and worst-case $\Omega(k\ell)$ for Perceptron
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos- const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update: $r_i^+(t) = \exp(-2\eta(\hat{y}(t) - y(t)) x_i(t)U)$ $r_i^-(t) = 1/r_i^+(t)$ $w_i^+(t+1) = U \cdot \frac{w_i^+(t)r_i^+(t)}{\sum_{j=1}^{\ell}(w_i^+(t)r_i^+(t) + w_i^-(t)r_i^-(t))}$ U and denominator normalize wts for proof of error bound	Winnow/Exponentiated Gradient Miscellany• Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron• Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries)• E.g. $x_i \in \{0, 1\}$, x's are labelled by a monotone k-disjunction over ℓ attributes, $k \ll \ell$ - Remaining $\ell - k$ are irrelevant- E.g. $x_5 \lor x_9 \lor x_{12}$, $\ell = 150$, $k = 3$ - For disjunctions, number of on-line prediction mistakes is $O(k \log \ell)$ for Winnow and worst-case $\Omega(k\ell)$ for Perceptron- So in worst case, need exponentially fewer updates for training in Win. than Percep.
Winnow/Exponentiated Gradient Negative Weights • Winnow and EG update wts by mult by a pos const, impossible to change sign – Weight vectors restricted to one quadrant • Solution: Maintain wt vectors $\mathbf{w}^+(t)$ and $\mathbf{w}^-(t)$ – Predict $\hat{y}(t) = (\mathbf{w}^+(t) - \mathbf{w}^-(t)) \cdot \mathbf{x}(t)$ – Update: $r_i^+(t) = \exp(-2\eta(\hat{y}(t) - y(t)) x_i(t)U)$ $r_i^-(t) = 1/r_i^+(t)$ $w_i^+(t+1) = U \cdot \frac{w_i^+(t)r_i^+(t)}{\sum_{j=1}^{\ell} (w_i^+(t)r_i^+(t) + w_i^-(t)r_i^-(t))}$ <i>U</i> and denominator normalize wts for proof of error bound Kivinen & Warmuth, "Additive Versus Exponen- tiated Gradient Updates for Linear Prediction." <i>Information and Computation</i> , 132(1):1–64, Jan. 1997. [see web page]	Winnow/Exponentiated Gradient Miscellany• Winnow and EG are muliplicative weight update schemes versus additive weight update schemes, e.g. Perceptron• Winnow and EG work well when most attributes (features) are irrelevant, i.e. optimal weight vector w* is sparse (many 0 entries)• E.g. $x_i \in \{0, 1\}$, x's are labelled by a monotone k-disjunction over ℓ attributes, $k \ll \ell$ - Remaining $\ell - k$ are irrelevant- E.g. $x_5 \lor x_9 \lor x_{12}$, $\ell = 150$, $k = 3$ - For disjunctions, number of on-line prediction mistakes is $O(k \log \ell)$ for Winnow and worst-case $\Omega(k\ell)$ for Perceptron- So in worst case, need exponentially fewer updates for training in Win. than Percep.• Other bounds exist for real-valued inputs and outputs

Non-Linearly Separable Classes

- What if no hyperplane completely separates the classes?
- Add extra inputs that are nonlinear combinations of original inputs (Section 4.14)
 - E.g. attribs. x_1 and x_2 , so try $\mathbf{x} = \begin{bmatrix} x_1, x_2, x_1x_2, x_1^2, x_2^2, x_1^2x_2, x_1x_2^2 \end{bmatrix}^T$
 - Perhaps classes linearly separable in new feature space
 - Useful esp. with Winnow/EG logarithmic bounds
 - Kernel functions/SVMs
- Pocket algorithm (p. 63) guarantees convergence to best hyperplane
- Winnow's & EG's agnostic results
- Least squares methods (Sec. 3.4)
- Networks of classifiers (Ch. 4)

13

Non-Linearly Separable Classes Least Squares Methods

• Recall from Slide 7:

$$w_i(t+1) = w_i(t) + \eta \left(y(t) - \sum_{i=1}^{\ell} w_i(t) x_i(t) \right) x_i(t)$$
$$= w_i(t) + \eta \left(y(t) - \mathbf{w}(t)^T \cdot \mathbf{x}(t) \right) x_i(t)$$

• If we <u>don't</u> threshold dot product during training and allow η to vary each trial (i.e. substitute η_t), get* Eq. 3.38, p. 69:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta_t \mathbf{x}(t) \left(y(t) - \mathbf{w}(t)^T \cdot \mathbf{x}(t) \right)$$

- This is Least Mean Squares (LMS) Algorithm
- If e.g. $\eta_t = 1/t$, then

$$\lim_{t \to \infty} P\left(\mathbf{w}(t) = \mathbf{w}^*\right) = \mathbf{1},$$

where

$$\mathbf{w}^{*} = \mathop{\mathrm{argmin}}_{\mathbf{w} \in \Re^{\ell}} \left\{ \mathsf{E} \left[\left| \boldsymbol{y} - \mathbf{w}^{T} \cdot \mathbf{x} \right|^{2} \right] \right\}$$

is vector minimizing mean square error (MSE)

*Note that here w(t) is weight <u>before</u> trial t. In book it is weight <u>after</u> trial t.

15

Non-Linearly Separable Classes

Winnow's Agnostic Results

• Winnow's total number of prediction mistakes loss (in <u>on-line setting</u>) provably not much worse than best linear classifier



- Loss bound related to performance of best classifier and total distance under $\|\cdot\|_1$ that feature vectors must be moved to make best classifier perfect [Littlestone, COLT '91]
- Similar bounds for EG [Kivinen & Warmuth]

14



$\mathbf{w} = \begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 & \mathbf{w}_3 \\ \mathbf{w} = \begin{bmatrix} \overline{w}_{11}, \overline{w}_{12}, \overline{w}_{10}, \overline{w}_{21}, \overline{w}_{22}, \overline{w}_{20}, \overline{w}_{31}, \overline{w}_{32}, \overline{w}_{30} \end{bmatrix}^T$ • Now if $\mathbf{w}^{*T} \cdot \mathbf{x}_1 > 0$ and $\mathbf{w}^{*T} \cdot \mathbf{x}_2 > 0$, then $\begin{pmatrix} \text{Class} & \text{Binary Encoding} \\ \text{Classifier 1} & \text{Classifier 2} \\ \hline{\omega_1} & 0 & 0 \\ \overline{\omega_2} & 0 & 1 \\ \overline{\omega_3} & 1 & 0 \\ \overline{\omega_4} & 1 & 1 \\ \end{pmatrix}$ $\begin{pmatrix} \ell+1 \\ \omega_3 \\ \omega_4 \\ \ell \\ 1 \\ \ell \\ \ell$	Multiclass learning Kessler's Construction (cont'd) • So map x to $x_1 = \begin{bmatrix} 0 & \text{neg} & \text{pad} \\ 2, 2, 1, -2, -2, -1, 0, 0, 0 \end{bmatrix}^T$ $x_2 = \begin{bmatrix} 2, 2, 1, 0, 0, 0, -2, -2, -1 \end{bmatrix}^T$ and let	 Multiclass learning Error-Correcting Output Codes (ECOC) Since Win. & Percep. learn binary functions, learn individual bits of <u>binary encoding</u> of classes E.g. M = 4, so use two linear classifiers:
 In general, map (ℓ+1)×1 feature vector x to x₁,x_{M-1}, each of size (ℓ+1)M×1 x ∈ ω_i ⇒ x in <i>i</i>th block and -x in <i>j</i>th block, (rest are 0s). Repeat for all <i>j</i> ≠ <i>i</i> Now train to find weights for new vector space 17 Problem: Sensitive to individual classifier energy of encodings per class to improve robustness Similar to principle of error-correcting output codes used in communication networks [Dietterich & Bakiri, 1995] General-purpose, independent of learner 	and let $\mathbf{w} = \begin{bmatrix} \mathbf{w}_{11}, \mathbf{w}_{12}, \mathbf{w}_{10}, \mathbf{w}_{21}, \mathbf{w}_{22}, \mathbf{w}_{20}, \mathbf{w}_{31}, \mathbf{w}_{32}, \mathbf{w}_{30} \end{bmatrix}^{T}$ • Now if $\mathbf{w}^{*T} \cdot \mathbf{x}_{1} > 0$ and $\mathbf{w}^{*T} \cdot \mathbf{x}_{2} > 0$, then $\begin{pmatrix} \ell+1 \\ \sum_{i=1}^{l} w_{1i}^{*} x_{i} > \sum_{i=1}^{\ell+1} w_{2i}^{*} x_{i} \text{AND} \sum_{i=1}^{\ell+1} w_{1i}^{*} x_{i} > \sum_{i=1}^{\ell+1} w_{3i}^{*} x_{i}$ • In general, map $(\ell + 1) \times 1$ feature vector \mathbf{x} to $\mathbf{x}_{1}, \dots, \mathbf{x}_{M-1}$, each of size $(\ell + 1)M \times 1$ • $\mathbf{x} \in \omega_{i} \Rightarrow \mathbf{x}$ in <i>i</i> th block and $-\mathbf{x}$ in <i>j</i> th block, (rest are 0s). Repeat for all $j \neq i$ • Now train to find weights for new vector space	ClassBinary Encoding Classifier 1 u_1 0 u_2 0 u_3 1 u_4 1and train simultaneously• Problem: Sensitive to individual classifier errors, so use a set of encodings per class to improve robustness• Similar to principle of error-correcting output codes used in communication networks [Dietterich & Bakiri, 1995]• General-purpose, independent of learner