

Linear Classifiers

Outline

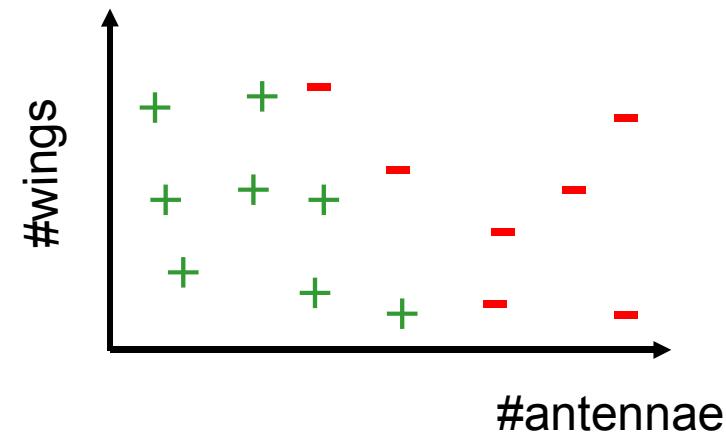
- Framework
- “Exact”
 - Minimize Mistakes (Perceptron Training)
 - Matrix inversion
- “Logistic Regression” Model
 - Max Likelihood Estimation (MLE) of $P(y | x)$
 - Gradient descent (MSE; MLE)
- “Linear Discriminant Analysis”
 - Max Likelihood Estimation (MLE) of $P(y, x)$
 - Direct Computation

Diagnosing Butterfly-itis



Classifier: Decision Boundaries

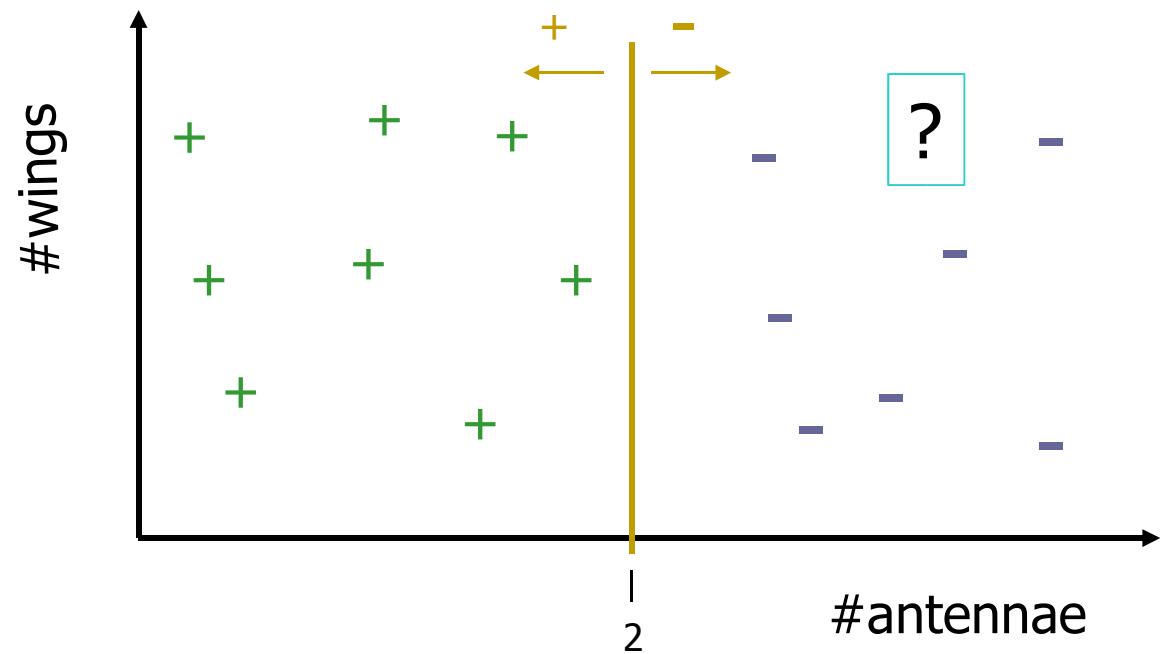
- **Classifier:** partitions input space X into “decision regions”



- *Linear threshold unit* has a linear decision boundary
- Defn: Set of points that can be separated by linear decision boundary is “linearly separable”

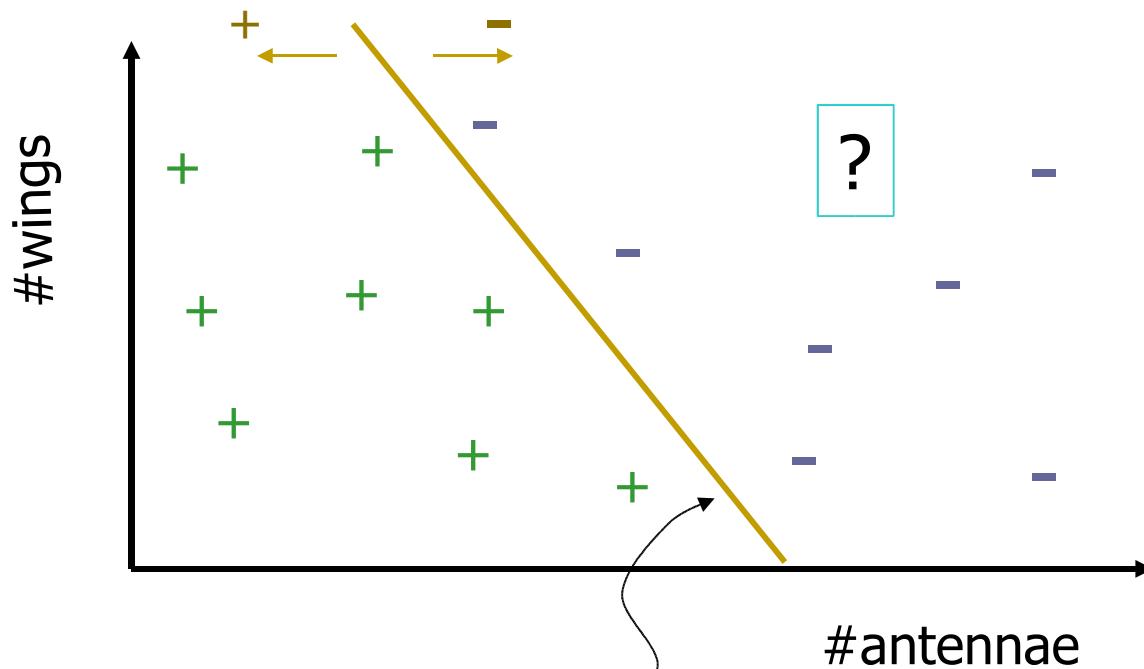
Linear Separators

- Draw “separating line”



- If $\# \text{antennae} \leq 2$, then butterfly-it is
- So ? is Not butterfly-it is.

Can be “angled”...



$$2.3 \times \#w + 7.5 \times \#a + 1.2 = 0$$

- If $2.3 \times \#Wings + 7.5 \times \#antennae + 1.2 > 0$
then butterfly-it-is

Linear Separators, in General

- Given data (many features)

F ₁	F ₂	...	F _n	Class
35	95	...	3	No
22	80	...	-2	Yes
:	:		:	:
10	50	...	1.9	No

- find “weights” $\{w_1, w_2, \dots, w_n, w_0\}$

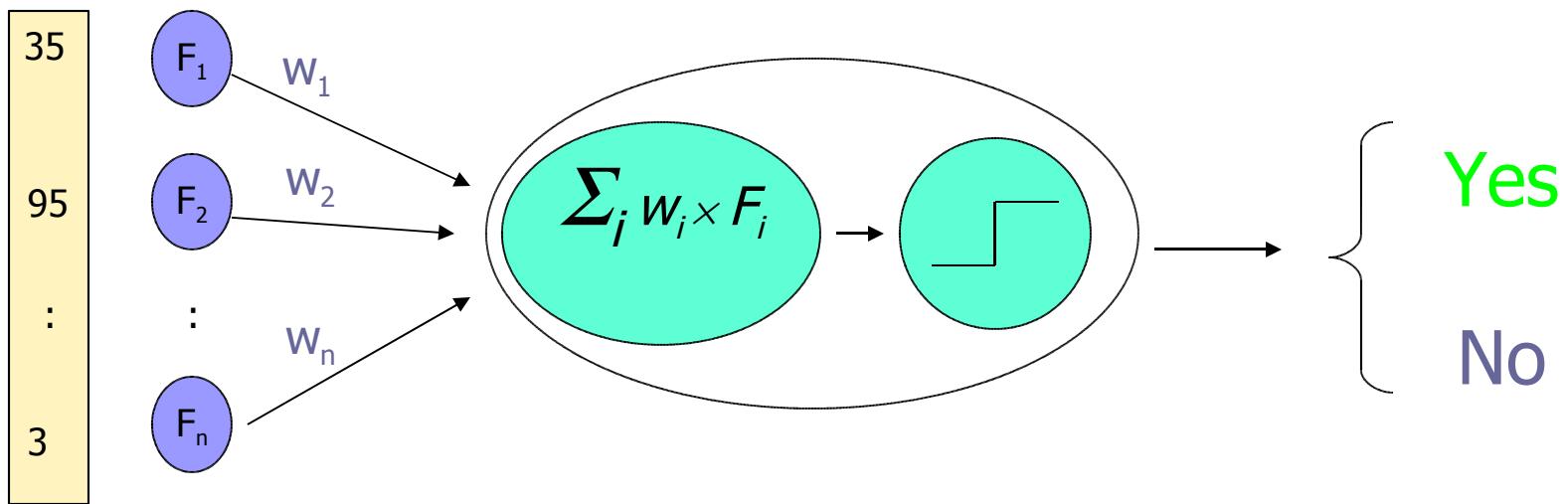
such that

$$w_1 \times F_1 + \dots + w_n \times F_n + w_0 > 0$$

means

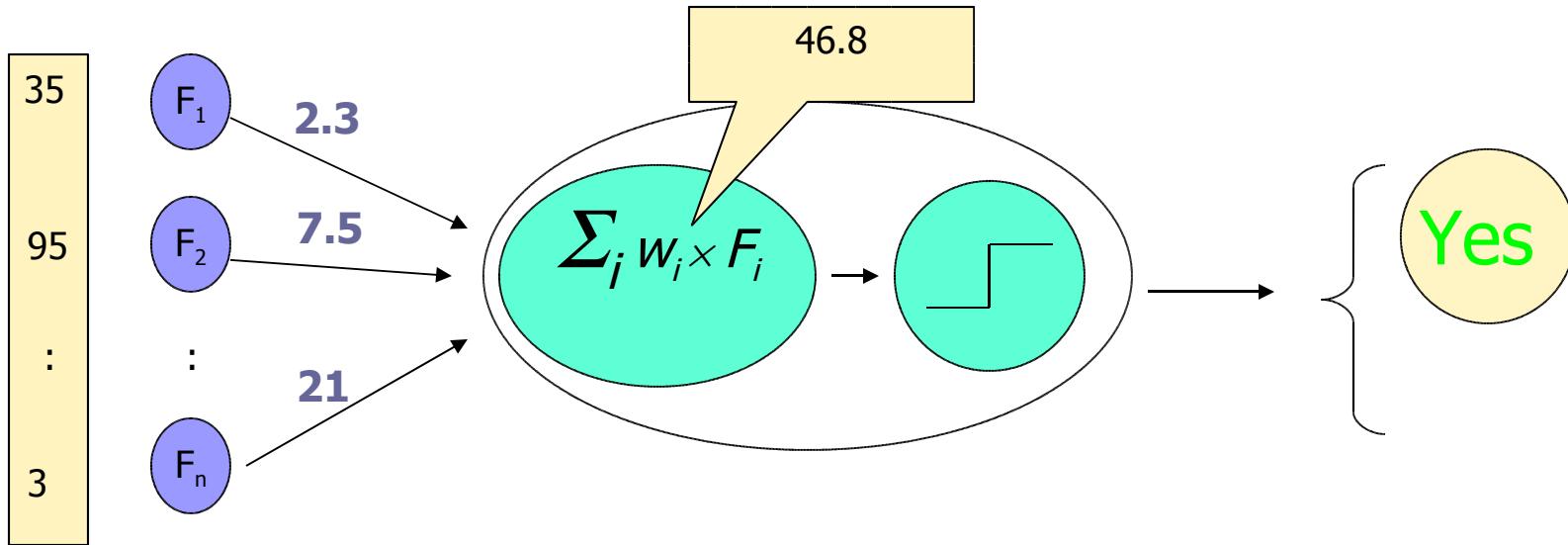
Class = Yes

Linear Separator



Just view $F_0 = 0$, so $w_0 \dots$

Linear Separator



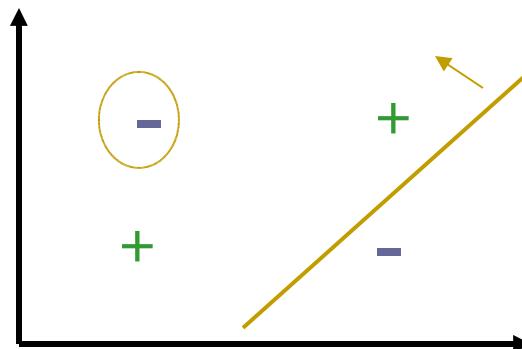
- Performance
 - Given $\{w_i\}$, and values for instance, compute response
- Learning
 - Given labeled data, find “correct” $\{w_i\}$
- Linear Threshold Unit ... “Perceptron”

Linear Separators – Facts

■ GOOD NEWS:

- If data is linearly separated,
- Then FAST ALGORITHM finds correct $\{w_i\}$!

■ But...

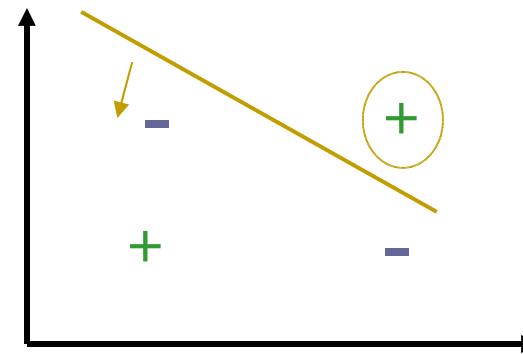


Linear Separators – Facts

■ GOOD NEWS:

- If data is linearly separated,
- Then FAST ALGORITHM finds correct $\{w_i\}$!

■ But...



- Some “data sets” are
NOT linearly separable!

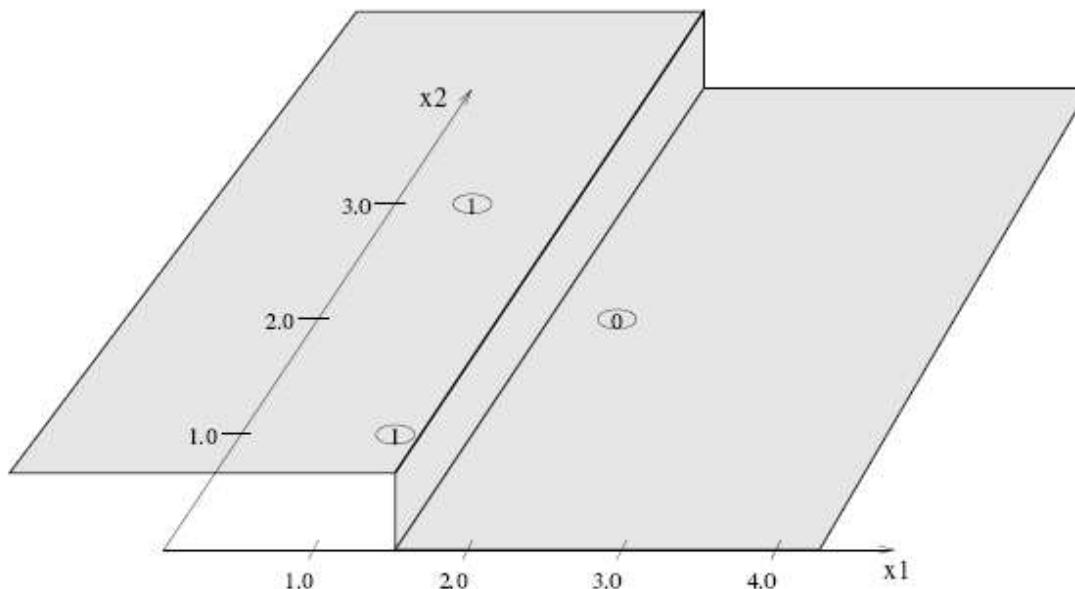
Stay tuned!

Geometric View

- Consider 3 training examples:

```
( [1.0, 1.0]; 1 )  
( [0.5; 3.0]; 1 )  
( [2.0; 2.0]; 0 )
```

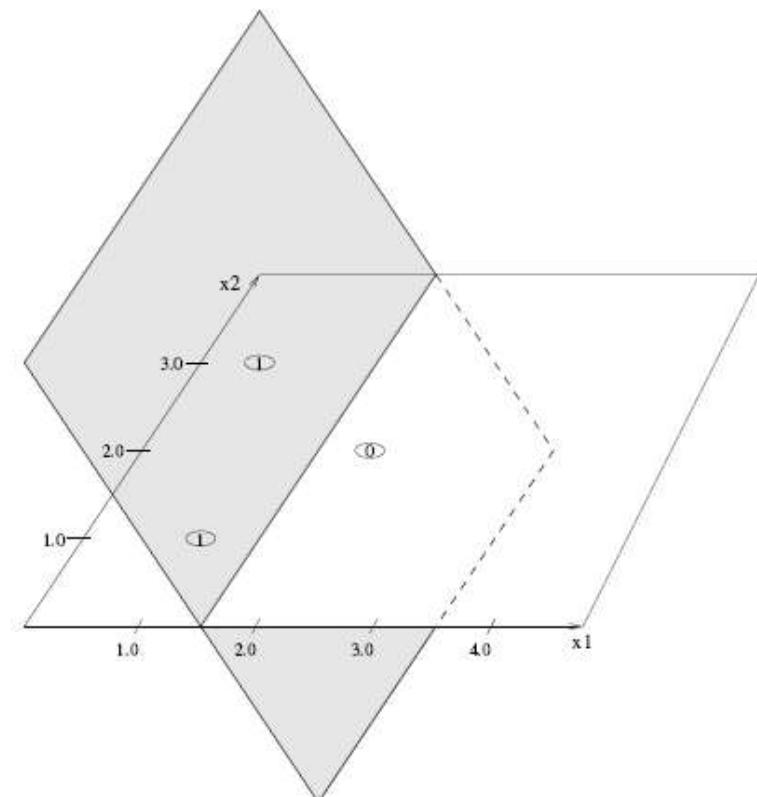
- Want classifier that looks like. . .



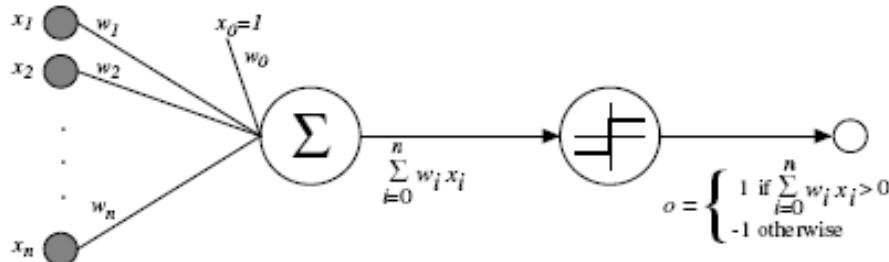
Linear Equation is Hyperplane

- Equation $\mathbf{w} \cdot \mathbf{x} = \sum_i w_i \cdot x_i$ is plane

$$y(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$



Linear Threshold Unit: “Perceptron”



$$\begin{aligned} o_w(x_1, \dots, x_n) &= \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases} \\ &= \text{sign}((w_0, w_1, \dots, w_n) \cdot (1, x_1, \dots, x_n)) \end{aligned}$$

- Squashing function:

$\text{sgn}: \mathbb{R} \rightarrow \{-1, +1\}$

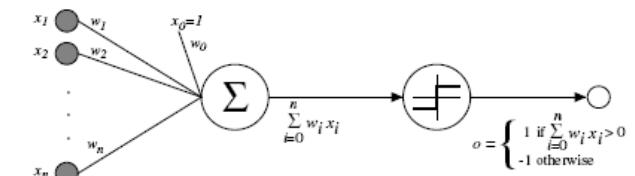
$$\text{sgn}(r) = \begin{cases} 1 & \text{if } r > 0 \\ 0 & \text{otherwise} \end{cases}$$

(“heaviside”)

- Actually $\mathbf{w} \cdot \mathbf{x} > b$ but . . .

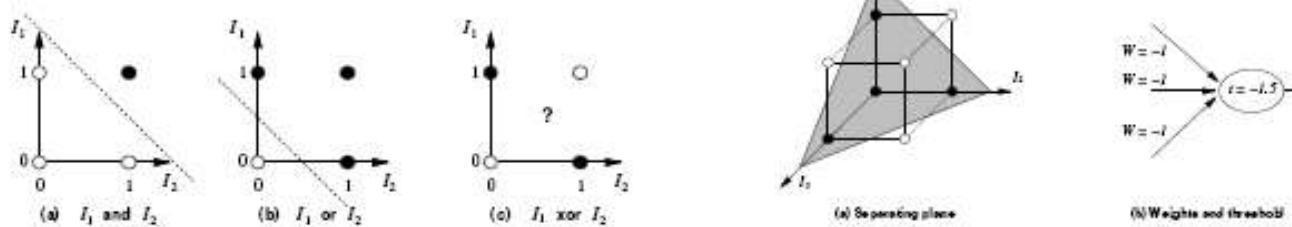
Create extra input x_0 fixed at 1

Corresponding w_0 corresponds to $-b$



Learning Perceptrons

- Can represent Linearly-Separated surface
... any hyper-plane between two half-spaces...



- Remarkable learning algorithm: [Rosenblatt 1960]

If function f can be represented by perceptron,
then \exists learning alg guaranteed to quickly converge to f !

⇒ enormous popularity, early / mid 60's

- But some simple fns cannot be represented
(Boolean XOR) [Minsky/Papert 1969]
- Killed the field temporarily!

Perceptron Learning

- Hypothesis space is . . .
 - **Fixed Size:**
 $\exists O(2^{n^2})$ distinct perceptrons over n boolean features
 - **Deterministic**
 - **Continuous** Parameters
- Learning algorithm:
 - Various: **Local** search, **Direct** computation, . . .
 - **Eager**
 - **Online / Batch**

Task

- Input: labeled data

	x_1	x_2	x_3	...	x_r	y
	16.1	-22.7	0.3	...	-4.0	+
	-7.3	-5.1	9.1	...	17.1	-
	:					:
	-16.2	-77.0	-1.2	...	-4.0	+

Transformed to

	x_0	x_1	x_2	x_3	...	x_r	y
	1	16.1	-22.7	0.3	...	-4.0	+
	1	-7.3	-5.1	9.1	...	17.1	-
	:	:					:
	1	-16.2	-77.0	-1.2	...	-4.0	+

- Output: $\mathbf{w} \in \Re^{r+1}$

Goal: Want \mathbf{w} s.t.

$$\forall i \text{ sgn}(\mathbf{w} \cdot [1, \mathbf{x}^{(i)}]) = y^{(i)}$$

□ . . . minimize mistakes wrt data . . .

Error Function

Given data $\{ [x^{(i)}, y^{(i)}] \}_{i=1..m}$, optimize...

- 1. Classification error

Perceptron Training; Matrix Inversion

$$err_{Class}(w) = \frac{1}{m} \sum_{i=1}^m I[y^{(i)} \neq o_w(x^{(i)})]$$

- 2. Mean-squared error

Matrix Inversion; Gradient Descent

$$err_{MSE}(w) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} [y^{(i)} - o_w(x^{(i)})]^2$$

- 3. (Log) Conditional Probability

MSE Gradient Descent; LCL Gradient Descent

$$LCL(w) = \frac{1}{m} \sum_{i=1}^m \log P_w(y^{(i)} | x^{(i)})$$

- 4. (Log) Joint Probability

Direct Computation

$$LL(w) = \frac{1}{m} \sum_{i=1}^m \log P_w(y^{(i)}, x^{(i)})$$

#1: Optimal Classification Error

- For each labeled instance $[x, y]$

$$\text{Err} = y - o_w(x)$$

$y = f(x)$ is target value

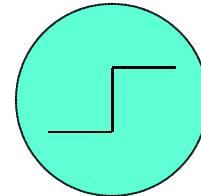
$o_w(x) = \text{sgn}(w \cdot x)$ is perceptron output

- Idea: Move weights in appropriate direction, to push $\text{Err} \rightarrow 0$
- If $\text{Err} > 0$ (error on POSITIVE example)
 - need to increase $\text{sgn}(w \cdot x)$
→ need to increase $w \cdot x$
 - Input j contributes $w_j \cdot x_j$ to $w \cdot x$
 - if $x_j > 0$, increase
 - if $x_j < 0$, decrease

$$\Rightarrow w_j \leftarrow w_j + x_j$$

If $\text{Err} < 0$ (error on NEGATIVE example)
 $\Rightarrow w_j \leftarrow w_j - x_j$

$$\Rightarrow w_j \leftarrow w_j + x_j$$



#1a: Mistake Bound Perceptron Alg

Initialize $w = 0$

Do until bored

Predict "+" iff $w \cdot x > 0$

else "-"

Mistake on positive: $w \leftarrow w + x$

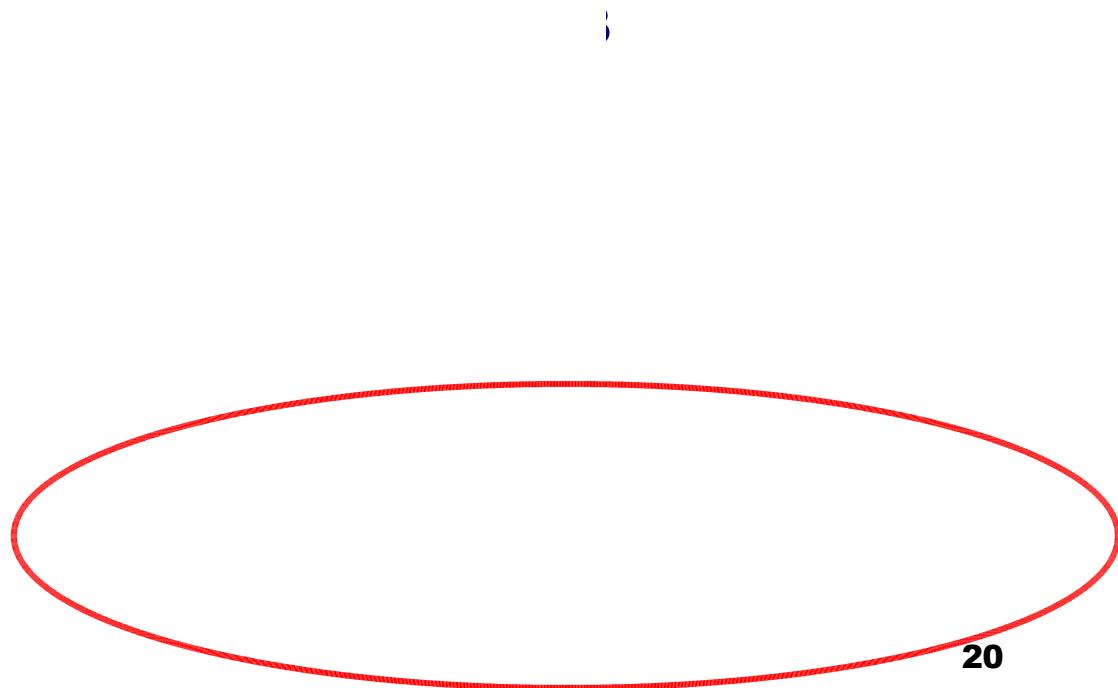
Mistake on negative: $w \leftarrow w - x$

Weights	Instance	Action
[0 0 0]		#1

Orig Data		
$\langle \langle x_1 \ x_2 \rangle \ c(x) \rangle$		
0	0	+
1	0	-
1	1	+

Data + " $x_0 = 1$ "

$\langle \langle x_0 \ x_1 \ x_2 \rangle \ c(x) \rangle$			
$i_1 :$	1	0	0
$i_2 :$	1	1	0
$i_3 :$	1	1	1

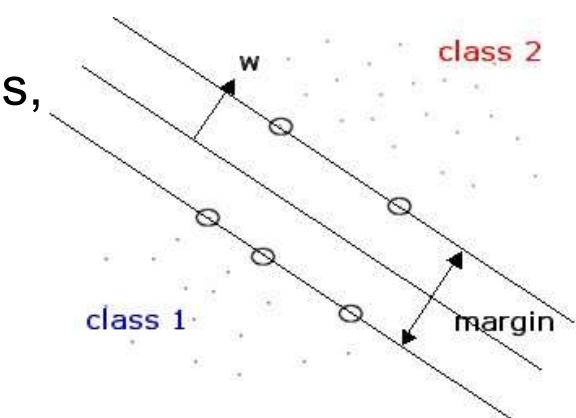


Mistake Bound Theorem

Theorem: [Rosenblatt 1960]

If data is consistent w/some linear threshold \mathbf{w} ,
then number of mistakes is $\leq (1/\Delta)^2$,
where $\Delta = \min_x \frac{|\mathbf{w} \cdot \mathbf{x}|}{|\mathbf{w}| \times |\mathbf{x}|}$

- Δ measures “wiggle room” available:
If $|\mathbf{x}| = 1$, then Δ is max, over all consistent planes,
of minimum distance of example to that plane
- \mathbf{w} is \perp to separator, as $\mathbf{w} \cdot \mathbf{x} = 0$ at boundary
- So $|\mathbf{w} \cdot \mathbf{x}|$ is projection of \mathbf{x} onto plane,
PERPENDICULAR to boundary line
... ie, is distance from \mathbf{x} to that line (once normalized)



Proof of Convergence

For simplicity:

0. Use $x_0 \equiv 1$, so target plane goes thru 0
1. Assume target plane doesn't hit any examples
2. Replace negative point $\langle (x_0, x_1, \dots, x_n) \ 0 \rangle$
by positive point $\langle (-x_0, -x_1, \dots, -x_n) \ 1 \rangle$
3. Normalize all examples to have length 1

\mathbf{x} wrong wrt \mathbf{w} iff $\mathbf{w} \cdot \mathbf{x} < 0$

- Let \mathbf{w}^* be unit vector rep'ning target plane

$$\Delta = \min_{\mathbf{x}} \{ \mathbf{w}^* \cdot \mathbf{x} \}$$

Let \mathbf{w} be hypothesis plane

- Consider:

$$\frac{(\mathbf{w} \cdot \mathbf{w}^*)}{|\mathbf{w}|}$$

- On each mistake, add \mathbf{x} to \mathbf{w}

$$\mathbf{w} = \sum_{\{\mathbf{x} \mid \mathbf{x} \cdot \mathbf{w} < 0\}} \mathbf{x}$$

Proof (con't)

If w is mistake...

Numerator increases by $x \cdot w^* \geq \Delta$

(denominator)² becomes

$$(w+x)^2 = w^2 + x^2 + 2(w \cdot x) < w^2 + 1$$

as $w \cdot x < 0$

As initially $w = \langle 0, \dots, 0 \rangle$.

After m mistakes,

numerator is $\geq m \times \Delta$

$$(\text{denominator})^2 \leq 0 + \underbrace{1 + \dots + 1}_m = m$$

so denominator $\leq \sqrt{m}$

$$\frac{(w \cdot w^*)}{|w|}$$

$$\Delta = \min_x \{ w^* \cdot x \}$$

$$W = \sum_{\{x \mid x \cdot w < 0\}} X$$

- As $(w \cdot w^*)/|w| = \cos(\text{angle between } w \text{ and } w^*)$
it must be ≤ 1 , so
numerator \leq denominator

$$\Rightarrow \Delta * m \leq \sqrt{m} \Rightarrow m \leq \frac{1}{\Delta^2}$$

#1b: Perceptron Training Rule

- For each labeled instance $[x, y]$

$$\text{Err}([x, y]) = y - o_w(x) \in \{-1, 0, +1\}$$

- If $\text{Err}([x, y]) = 0$ **Correct!** ... Do nothing!

$$\Delta w = 0 \equiv \text{Err}([x, y]) \cdot x$$

- If $\text{Err}([x, y]) = +1$ **Mistake on positive!** Increment by $+x$

$$\Delta w = +x \equiv \text{Err}([x, y]) \cdot x$$

- If $\text{Err}([x, y]) = -1$ **Mistake on negative!** Increment by $-x$

$$\Delta w = -x \equiv \text{Err}([x, y]) \cdot x$$

In all cases... $\Delta w^{(i)} = \text{Err}([x^{(i)}, y^{(i)}]) \cdot x^{(i)} = [y^{(i)} - o_w(x^{(i)})] \cdot x^{(i)}$

- Batch

$$\Delta w_j = \sum_i \Delta w_j^{(i)}$$

$$= \sum_i x_j^{(i)} (y^{(i)} - o_w(x^{(i)}))$$

$$w_j += \eta \Delta w_j$$

η is **learning rate** (small pos “constant” ... $\approx 0.05?$)

feature j

$\mathbf{X}^{(i)} \rightarrow$

$\Delta\mathbf{w} \rightarrow$

		$\mathbf{X}^{(i)}_j$		

0. New \mathbf{w}
1. For each row i , compute
- $\Delta\mathbf{w} = 0$
 - $E^{(i)} = y^{(i)} - o_{\mathbf{w}}(\mathbf{x}^{(i)})$
 - $\Delta\mathbf{w} += E^{(i)} \mathbf{x}^{(i)}$
[... $\Delta w_j += E^{(i)} x_j^{(i)} \dots$]
2. Increment $\mathbf{w} += \eta \Delta\mathbf{w}$

$E^{(i)}$	

Correctness

- Rule is intuitive: **Climbs in correct direction. . .**
- Thrm: Converges to correct answer, if . . .
 - training data is linearly separable
 - sufficiently small η
- Proof: Weight space has **EXACTLY 1 minimum!**
(no non-global minima)
⇒ with enough examples, finds correct function!
- Explains early popularity
- If η too large, may overshoot
If η too small, takes too long
- So often $\eta = \eta(k)$. . . which decays with # of iterations, k ²⁶

#1c: Matrix Version?

Task: Given $\{ \langle \mathbf{x}^i, y^i \rangle \}$

$y^i \in \{-1, +1\}$ is label

Find w_i s.t.

$$\begin{aligned}y^1 &= w_0 + w_1 x_1^1 + \cdots + w_n x_n^1 \\y^2 &= w_0 + w_1 x_1^2 + \cdots + w_n x_n^2 \\&\vdots \\y^m &= w_0 + w_1 x_1^m + \cdots + w_n x_n^m\end{aligned}$$

- **Linear Equalities** $\mathbf{y} = \mathbf{Xw}$!

$$\mathbf{y} = [y^1, \dots, y^m]^\top$$

$$\mathbf{X} = \begin{pmatrix} 1 & x_1^1 & \cdots & x_n^1 \\ 1 & x_1^2 & \cdots & x_n^2 \\ \vdots & \vdots & & \vdots \\ 1 & x_1^m & \cdots & x_n^m \end{pmatrix}$$

$$\mathbf{w} = [w_0, w_1, \dots, w_n]^\top$$

- Solution: $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$!

Issues

Task: Given $\{ \langle \mathbf{x}^i, y^i \rangle \}$ $y^i \in \{-1, +1\}$ is label
Find w_i s.t.

$$\begin{aligned}y^1 &= w_0 + w_1 x_1^1 + \cdots + w_n x_n^1 \\y^2 &= w_0 + w_1 x_1^2 + \cdots + w_n x_n^2 \\\vdots \\y^m &= w_0 + w_1 x_1^m + \cdots + w_n x_n^m\end{aligned}$$

1. Why restrict to only $y^i \in \{-1, +1\}$?

- If from discrete set $y^i \in \{0, 1, \dots, m\}$:
General (non-binary) classification
- If ARBITRARY $y^i \in \Re$: Regression

2. What if NO \mathbf{w} works?

... \mathbf{X} is singular; overconstrained ...

Could try to minimize residual

NP-Hard!

$$\sum_i \mathbb{I}[y^{(i)} \neq \mathbf{w} \cdot \mathbf{x}^{(i)}]$$

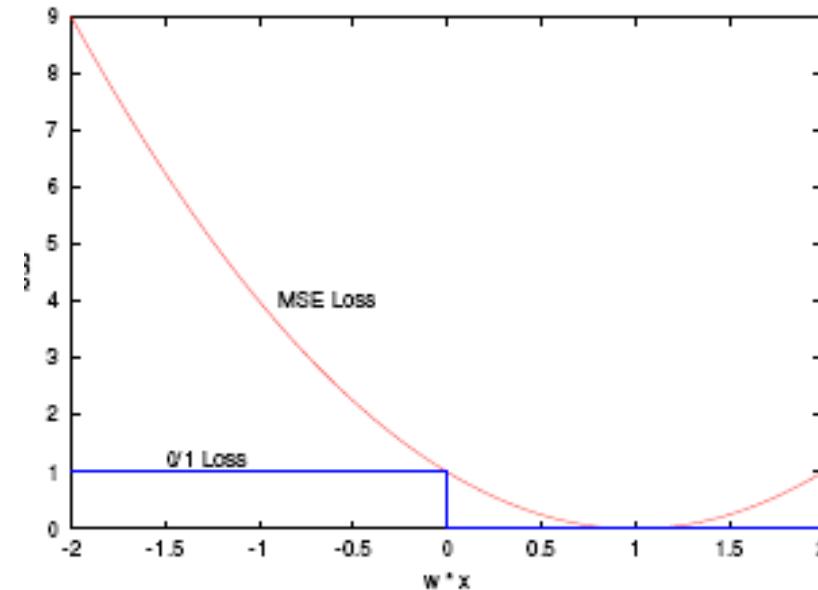
$$\| \mathbf{y} - \mathbf{X} \mathbf{w} \|_1 = \sum_i |y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)}|$$

$$\| \mathbf{y} - \mathbf{X} \mathbf{w} \|_2 = \sum_i (y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)})^2$$

Easy!

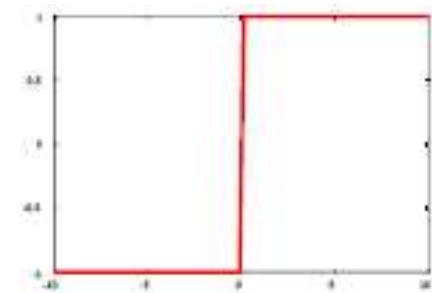
L_2 error vs 0/1-Loss

- “0/1 Loss function” not smooth, differentiable
- MSE error is smooth, differentiable... and is overbound...



Gradient Descent for Perceptron?

- Why not Gradient Descent for THRESHOLDed perceptron?
- Needs gradient (derivative), not



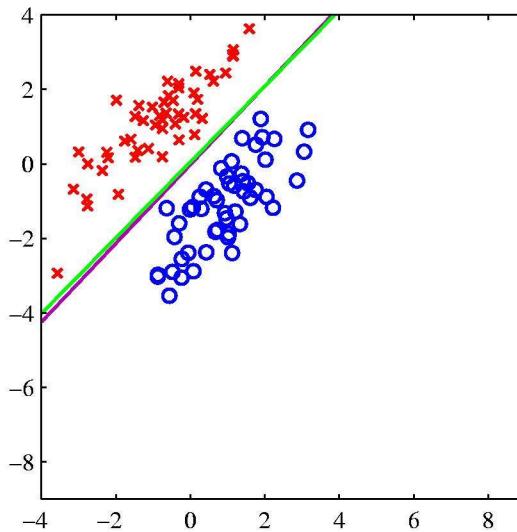
- Gradient Descent is General approach.
Requires
 - + continuously parameterized hypothesis
 - + error must be differentiable wrt parameters
- But...
 - can be slow (many iterations)
 - may only find LOCAL opt

#1. LMS version of Classifier

- View as Regression
 - Find “best” linear mapping \mathbf{w} from \mathbf{X} to \mathbf{Y}
 - $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \text{Err}_{\text{LMS}}(\mathbf{X}, \mathbf{Y})(\mathbf{w})$
 - $\text{Err}_{\text{LMS}}(\mathbf{X}, \mathbf{Y})(\mathbf{w}) = \sum_i (y^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)})^2$
 - Threshold: if $\mathbf{w} \cdot \mathbf{x} > 0.5$,
return 1;
else 0
- See Chapter 3

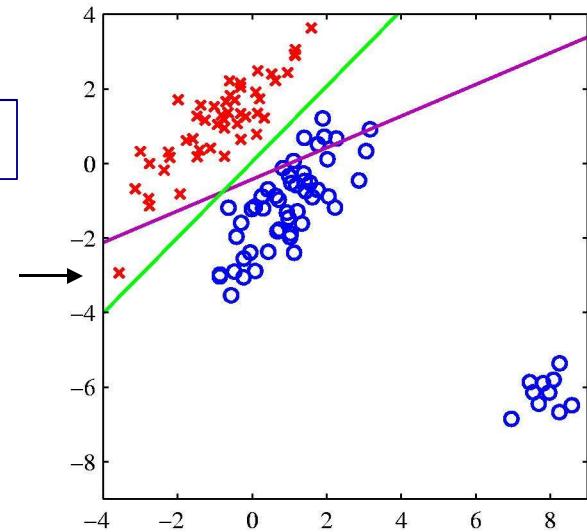
Use Linear Regression for Classification?

- 1. Use regression to find weights \mathbf{w}
 - 2. Classify new instance \mathbf{x} as $\text{sgn}(\mathbf{w} \cdot \mathbf{x})$
- But ... regression minimizes sum of squared errors on target function
... which gives strong influence to outliers



Great separation

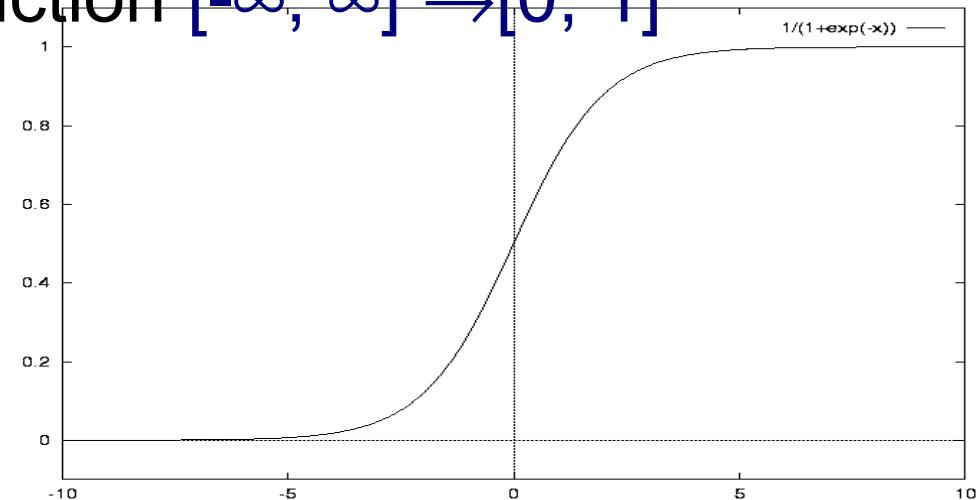
Bad separation



#3: Logistic Regression

- Want to compute $P_w(y=1 | \mathbf{x})$
... based on parameters \mathbf{w}
- But ...
 - $\mathbf{w} \cdot \mathbf{x}$ has range $[-\infty, \infty]$
 - probability must be in range $\in [0; 1]$
- Need “squashing” function $[-\infty, \infty] \rightarrow [0, 1]$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



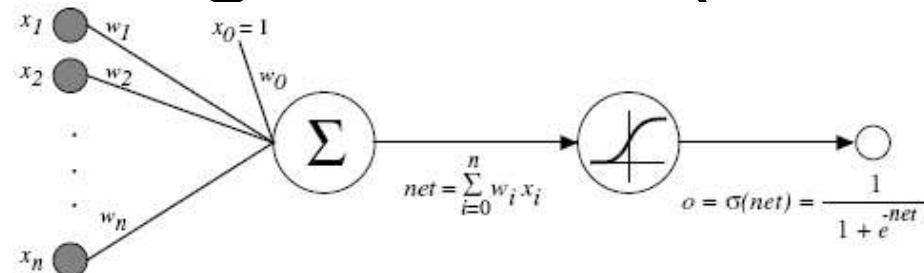
Alternative Derivation...

$$P(+y|x) = \frac{P(x|+y)P(+y)}{P(x|+y)P(+y) + P(x|-y)P(-y)}$$

$$= \frac{1}{1 + \exp(-a)}$$

$$a = \ln \frac{P(x|+y)P(+y)}{P(x|-y)P(-y)}$$

Logistic Regression (con't)



- Assume 2 classes:

$$P_w(+y|x) = \sigma(w \cdot x) = \frac{1}{1 + e^{-(x \cdot w)}}$$

$$P_w(-y|x) = 1 - \frac{1}{1 + e^{-(x \cdot w)}} = \frac{e^{-(x \cdot w)}}{1 + e^{-(x \cdot w)}}$$

- Log Odds:

$$\log \frac{P_w(+y|x)}{P_w(-y|x)} = x \cdot w$$

How to learn parameters w ?

- ... depends on goal?

- A: Minimize MSE?

$$\sum_i (y^{(i)} - o_w(\mathbf{x}^{(i)}))^2$$

- B: Maximize likelihood?

$$\sum_i \log P_w(y^{(i)} | \mathbf{x}^{(i)})$$

MSError Gradient for Sigmoid Unit

- Error: $\sum_j (y^{(j)} - o_w(\mathbf{x}^{(j)}))^2 = \sum_j E^{(j)}$

For single training instance

- Input:** $\mathbf{x}^{(j)} = [x_1^{(j)}, \dots, x_k^{(j)}]$
- Computed Output:** $o^{(j)} = \sigma(\sum_i x_i^{(j)} \cdot w_i) = \sigma(z^{(j)})$

□ where $z^{(j)} = \sum_i x_i^{(j)} \cdot w_i$ using current $\{w_i\}$

- Correct output:** $y^{(j)}$

Stochastic Error Gradient ..

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \left[\frac{1}{2} (o - y)^2 \right] = \frac{1}{2} \left[2(o - y) \frac{\partial}{\partial w_i} (o - y) \right] \\ &= (o - y) \left(\frac{\partial o}{\partial w_i} \right) = (o - y) \frac{\partial \sigma(z)}{\partial z} \frac{\partial z}{\partial w_i}\end{aligned}$$

Derivative of Sigmoid

$$\begin{aligned}\frac{d}{da} \sigma(a) &= \frac{d}{da} \frac{1}{(1+e^{-a})} \\&= \frac{-1}{(1+e^{-a})^2} \frac{d}{da}(1+e^{-a}) = \frac{-1}{(1+e^{-a})^2}(-e^{-a}) \\&= \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{(1+e^{-a})} \frac{e^{-a}}{(1+e^{-a})} = \sigma(a) [1 - \sigma(a)]\end{aligned}$$

Updating LR Weights (MSE)

- $\frac{\partial E}{\partial w_i} = (o - y) \frac{\partial \sigma(z)}{\partial z} \frac{\partial z}{\partial w_i}$

- Using:

$$\begin{aligned}\frac{\partial \sigma(z)}{\partial z} &= \sigma(z)(1 - \sigma(z)) = o(1 - o) \\ \frac{\partial z}{\partial w_i} &= \frac{\partial(\sum_i w_i \cdot x_i)}{\partial w_i} = x_i\end{aligned}$$

$$\Rightarrow \boxed{\frac{\partial E^{(j)}}{\partial w_i} = (o^{(j)} - y^{(j)}) o^{(j)} (1 - o^{(j)}) x_i^{(j)}}$$

Note: As already computed $o^{(j)} = \sigma(z^{(j)})$ to get answer,
trivial to compute $\sigma'(z^{(j)}) = \sigma(z^{(j)})(1 - \sigma(z^{(j)}))$!

- Update $w_i += \Delta w_i$ where

$$\boxed{\Delta w_i = \eta \cdot \frac{\partial E^{(j)}}{\partial w_i}}$$

B: Or... Learn Conditional Probability

- As fitting probability distribution,
better to return probability distribution ($\approx \mathbf{w}$)
that is **most likely, given training data, S**

$$\begin{aligned}\text{Goal: } \mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w} | S) \\ &= \operatorname{argmax}_{\mathbf{w}} \frac{P(S | \mathbf{w})P(\mathbf{w})}{P(S)} \\ &= \operatorname{argmax}_{\mathbf{w}} P(S | \mathbf{w})P(\mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} P(S | \mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \log P(S | \mathbf{w})\end{aligned}$$

Bayes Rules

As $P(S)$ does not depend on \mathbf{w}

As $P(\mathbf{w})$ is uniform

As \log is monotonic

ML Estimation

- $P(S | w)$ \equiv likelihood function

$$L(w) = \log P(S | w)$$

- $w^* = \operatorname{argmax}_w L(w)$

is “maximum likelihood estimator” (MLE)

Computing the Likelihood

- As training examples $[\mathbf{x}^{(i)}, y^{(i)}]$ are iid
 - drawn independently from same (unknown) prob $P_w(\mathbf{x}, y)$
- $\log P(S | \mathbf{w}) = \log \prod_i P_w(\mathbf{x}^{(i)}, y^{(i)})$ $= \sum_i \log P_w(\mathbf{x}^{(i)}, y^{(i)})$ $= \sum_i \log P_w(y^{(i)} | \mathbf{x}^{(i)}) + \sum_i \log P_w(\mathbf{x}^{(i)})$
- Here $P_w(\mathbf{x}^{(i)}) = 1/n \dots$
not dependent on \mathbf{w} , over empirical sample S
- $\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \sum_i \log P_w(y^{(i)} | \mathbf{x}^{(i)})$

Fit Logistic Regression... by Gradient Ascent

- Want $\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} J(\mathbf{w})$

- $J(\mathbf{w}) = \sum_i r(y^{(i)}, \mathbf{x}^{(i)}, \mathbf{w})$

- For $y \in \{0, 1\}$

$$r(y, \mathbf{x}, \mathbf{w}) = \log P_{\mathbf{w}}(y | \mathbf{x}) =$$

$$y \log(P_{\mathbf{w}}(y=1 | \mathbf{x})) + (1 - y) \log(1 - P_{\mathbf{w}}(y=1 | \mathbf{x}))$$

$$\frac{\partial J(w)}{\partial w_j} = \sum_i \frac{\partial r(y^{(i)}, x^{(i)}, w)}{\partial w_j}$$

- So climb along...

Gradient Descent ...

$$\begin{aligned}\frac{\partial r(y, x, w)}{\partial w_j} &= \frac{\partial}{\partial w_j} [y \log(p_1) + (1-y)\log(1-p_1)] \\&= \frac{y}{p_1} \frac{\partial p_1}{\partial w_j} + (-1) \times \frac{1-y}{1-p_1} \frac{\partial p_1}{\partial w_j} = \frac{y-p_1}{p_1(1-p_1)} \frac{\partial p_1}{\partial w_j}\end{aligned}$$

$$\begin{aligned}\frac{\partial p_1}{\partial w_j} &= \frac{\partial P_w(y=1|x)}{\partial w_j} = \frac{\partial}{\partial w_j} (\sigma(x \cdot w)) \\&= \sigma(x \cdot w) [1 - \sigma(x \cdot w)] \frac{\partial}{\partial w_j} (x \cdot w) = p_1(1-p_1) \cdot x_j^{(i)}\end{aligned}$$

$$\begin{aligned}\frac{\partial J(w)}{\partial w_j} &= \sum_i \frac{\partial r(y^{(i)}, x^{(i)}, w)}{\partial w_j} = \sum_i \frac{y^{(i)} - p_1}{p_1(1-p_1)} p_1(1-p_1) \cdot x_j^{(i)} \\&= \sum_i (y^{(i)} - P_w(y=1|x)) \cdot x_j^{(i)}\end{aligned}$$

Gradient Ascent for Logistic Regression (MLE)

Given: training examples $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$, $i = 1..N$

Set initial weight vector $\mathbf{w} = \langle 0, 0, 0, 0, \dots, 0 \rangle$

Repeat until convergence

Let gradient vector $\Delta\mathbf{w} = \langle 0, 0, 0, 0, \dots, 0 \rangle$

For $i = 1$ to N **do**

$$p_1^{(i)} = 1/(1 + \exp[\mathbf{w} \cdot \mathbf{x}^{(i)}])$$

$$\text{error}_i = y^{(i)} - p_1^{(i)}$$

For $j = 1$ to n **do**

$$\Delta\mathbf{w}_j += \text{error}_i \cdot x_{ij}$$

$\mathbf{w} += \eta \Delta\mathbf{w}$ % step in direction of increasing gradient

Comments on MLE Algorithm

- This is BATCH;
 \exists obvious online alg
(stochastic gradient ascent)
- Can use second-order (Newton-Raphson) alg for faster convergence
 - weighted least squares computation;
aka
“Iteratively-Reweighted Least Squares” (IRLS)

Use Logistic Regression for Classification

- Return YES iff

$$P(y=1|x)$$

$$\dot{P}(y=0|x)$$

$$\frac{P(y=1|x)}{P(y=0|x)}$$

$$\dot{1}$$

$$\ln \frac{P(y=1|x)}{P(y=0|x)}$$

$$\dot{0}$$

$$\ln \frac{1/(1+\exp(-w \cdot x))}{\exp(-w \cdot x)/(1+\exp(-w \cdot x))}$$

$$\dot{0}$$

$$\ln \frac{1}{\exp(-w \cdot x)} = w \cdot x > 0$$

Logistic Regression learns a LTU!

Logistic Regression for $K > 2$ Classes

- To handle $K > 2$ classes
 - Let class K be “reference”
 - Represent each other class $k \neq K$ as logistic function of odds of class k versus class K :

$$\begin{aligned} \log \frac{P(y=1|x)}{P(y=K|x)} &= w_1 \cdot x \\ \log \frac{P(y=2|x)}{P(y=K|x)} &= w_2 \cdot x \\ &\vdots \\ \log \frac{P(y=K-1|x)}{P(y=K|x)} &= w_{K-1} \cdot x \end{aligned}$$

- Apply gradient ascent to learn all w_k weight vectors, in parallel.

- Conditional probabilities:

$$P(y=k|x) = \frac{\exp(w_k \cdot x)}{1 + \sum_{\ell=1}^{K-1} \exp(w_\ell \cdot x)}$$

and

$$P(y=K|x) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(w_\ell \cdot x)}$$

Note: $K-1$ different w_i weights,
... each of dimension $|x|$

Learning LR Weights

Task: Given data $\langle \langle \mathbf{x}^{(i)}, y^{(i)} \rangle \rangle$

$$\text{find } \mathbf{w} \text{ in } p_{\mathbf{w}}(y|\mathbf{x}) = \begin{cases} \frac{1}{1+\exp(-\mathbf{w} \cdot \mathbf{x})} & \text{if } y=1 \\ \frac{\exp(-\mathbf{w} \cdot \mathbf{x})}{1+\exp(-\mathbf{w} \cdot \mathbf{x})} & \text{if } y=0 \end{cases}$$

$$\text{s.t. } p_{\mathbf{w}}(y^{(i)}|\mathbf{x}^{(i)}) > \frac{1}{2} \quad \text{iff} \quad y^{(i)} = 1$$

Approach 1: MSE – “Neural nets”

$$\text{Minimize } \sum_i (o^{(i)} - y^{(i)})^2$$

Gradient:

$$\Delta \mathbf{w}^{(i)}_j = (o^{(i)} - y^{(i)}) o^{(i)} (1 - o^{(i)})$$

Approach 2: MLE – “Logistic Regression”

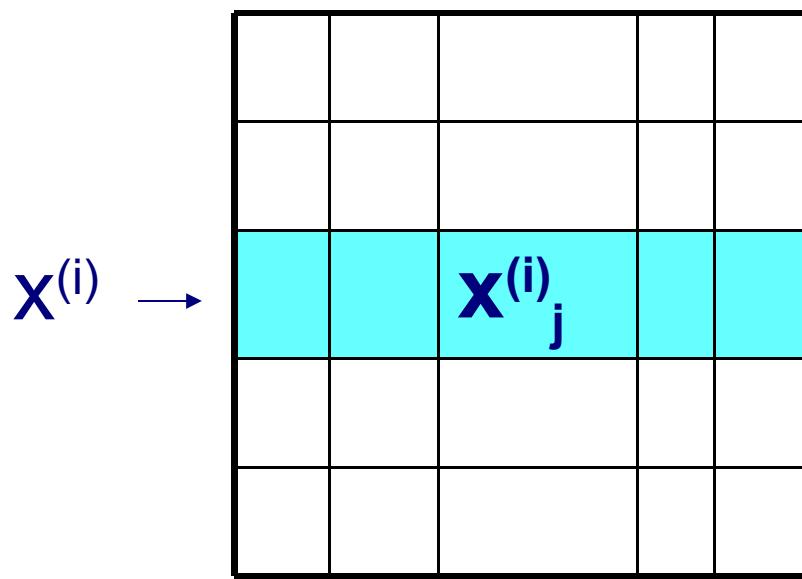
$$\text{Maximize } \sum_i p_{\mathbf{w}}(y|\mathbf{x})$$

Gradient:

$$\Delta \mathbf{w}^{(i)}_j = (y^{(i)} - p(1|\mathbf{x}^{(i)})) \mathbf{x}^{(i)}_j$$

(MaxProb)

feature j



0. New w
1. For each row i , compute
 - a. $\Delta w = 0$
 - b. $E^{(i)} = y^{(i)} - p(1|x^{(i)})$
 - c. $\Delta w += E^{(i)} X^{(i)}$
[... $\Delta w_j += E^{(i)} X^{(i)}_j$...]
2. Increment $w += \eta \Delta w$

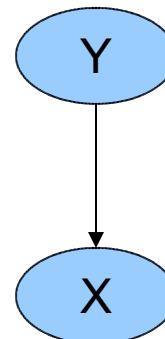
Logistic Regression Algs for LTUs

- Learns Conditional Probability Distribution $P(y | x)$
- **Local Search:**
Begin with initial weight vector;
iteratively modify to maximize objective function
log likelihood of the data
(ie, seek w s.t. probability distribution $P_w(y | x)$ is
most likely given data.)
- **Eager:** Classifier constructed from training examples,
which can then be discarded.
- **Online or batch**

#4: Linear Discriminant Analysis

- LDA learns joint distribution $P(y, \mathbf{x})$
 - As $P(y, \mathbf{x}) \neq P(y | \mathbf{x})$;
optimizing $P(y, \mathbf{x}) \neq$ optimizing $P(y | \mathbf{x})$
- “generative model”
 - $P(y, \mathbf{x})$ model of how data is generated
 - Eg, factor
 - $P(y, \mathbf{x}) = P(y) P(\mathbf{x} | y)$
 - $P(y)$ generates value for y ; then
 - $P(\mathbf{x} | y)$ generates value for \mathbf{x} given this y

- Belief net:



Linear Discriminant Analysis, con't

- $P(y, \mathbf{x}) = P(y)P(\mathbf{x} | y)$
- $P(y)$ is a simple discrete distribution
 - Eg: $P(y = 0) = 0.31; P(y = 1) = 0.69$
(31% negative examples; 69% positive examples)
- Assume $P(\mathbf{x} | y)$ is multivariate normal,
with mean μ_k and covariance Σ

$$P(\mathbf{x} | y = k) =$$

$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_k)^\top \Sigma^{-1} (\mathbf{x} - \mu_k) \right]$$

Estimating LDA Model

- Linear discriminant analysis assumes form

$$P(x, y) = P(y) \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_y)^\top \Sigma^{-1} (x - \mu_y) \right]$$

- μ_y is mean for examples belonging to class y ;
covariance matrix Σ is shared by all classes !
- Can estimate LDA directly:
 m_k = #training examples in class $y = k$

- Estimate of $P(y = k)$: $p_k = m_k / m$

$$\hat{\mu}_k = \frac{1}{m} \sum_{\{i: y_i = k\}} x_i$$

$$\hat{\Sigma} = \frac{1}{m} \sum_i (x_i - \hat{\mu}_{y_i})(x_i - \hat{\mu}_{y_i})^T$$

$$\hat{\mu}_{y_i}$$

(Subtract each x_i from corresponding $\hat{\mu}_{y_i}$ before taking outer product)

Example of Estimation

x_1	x_2	x_3	y
13.1	20.2	0.4	+
6.0	17.7	-4.2	+
8.2	18.2	-2.5	+
0.4	10.1	19.2	-
-4.2	12.8	5.1	-
-4.3	15.0	21.7	-
0.9	10.1	19.2	-

Note: do NOT pre-pend $x_0 = 1$!

- m=7 examples;
 $m_+ = 3$ positive; $m_- = 4$ negative
 $\Rightarrow p_+ = 3/7 \quad p_- = 4/7$
- Compute $\hat{\mu}_i$ over each class

$$\begin{aligned}\hat{\mu}_+ &= \frac{1}{3} \sum_{i: \langle y^{(i)} = + \rangle} \mathbf{x}^{(i)} \\ &= \frac{1}{3} \left([13.1, 20.2, 0.4] + [6.0, 17.7, -4.2] + [8.2, 18.2, -2.5] \right) \\ &= [9.1, 18.7, -2.1] \\ \hat{\mu}_- &= \frac{1}{4} \sum_{i: \langle y^{(i)} = - \rangle} \mathbf{x}^{(i)} = [-1.8, 12.0, 16.3]\end{aligned}$$

Estimation...

- Compute common $\hat{\Sigma}$

– “Normalize” each $z := x - \mu_y(x)$

$$\begin{aligned} z^{(1)} &:= [13.1, 20.2, 0.4] - [9.1, 18.7, -2.1] \\ &= [4.0, 1.5, -1.7] \end{aligned}$$

...

$$\begin{aligned} z^{(4)} &:= [0.4, 10.1, 19.2] - [-1.8, 12.0, 16.3] \\ &= [2.2, -1.9, 2.9] \end{aligned}$$

... $z^{(7)} := \dots$

– Compute covariance matrix, for each i :

For $x^{(1)}$, via $z^{(1)}$:

$$z^{(1)} \times z^{(1)\top} = \begin{bmatrix} 4.0 \\ 0.5 \\ -1.7 \end{bmatrix} \cdot [4.0, 0.5, -1.7]$$

$$= \begin{bmatrix} 4.0 \cdot 4.0 & 4.0 \cdot 0.5 & 4.0 \cdot -1.7 \\ 0.5 \cdot 4.0 & 0.5 \cdot 0.5 & 0.5 \cdot -1.7 \\ -1.7 \cdot 4.0 & -1.7 \cdot 0.5 & -1.7 \cdot -1.7 \end{bmatrix}$$

$$= \begin{bmatrix} 16.0 & 2.0 & -6.8 \\ 2.0 & 0.25 & -0.85 \\ -6.8 & -0.85 & -2.89 \end{bmatrix}$$

– Set $\hat{\Sigma} = \frac{1}{m} \sum_i z^{(i)} z^{(i)\top}$

x_1	x_2	x_3	y
13.1	20.2	0.4	+
6.0	17.7	-4.2	+
8.2	18.2	-2.5	+
0.4	10.1	19.2	-
-4.2	12.8	5.1	-
-4.3	15.0	21.7	-
0.9	10.1	19.2	-

Classifying, Using LDA

■ How to classify new instance, given estimates

Eg, $\hat{p}_+ = 3/7 \quad \hat{p}_- = 4/7$

$\{\hat{p}_i\}$ $\{\hat{\mu}_i\}$ $\hat{\Sigma}$

$$\begin{aligned} \star \hat{\mu}_+ &= [9.1, 18.7, -2.1] \\ \hat{\mu}_- &= [-1.8, 12.0, 16.3] \end{aligned}$$

$$\star \hat{\Sigma} = \begin{bmatrix} 7.22 & -1.31 & 6.35 \\ -1.31 & 2.91 & 0.32 \\ 6.35 & 0.32 & 26.03 \end{bmatrix}$$

■ Class for instance $\mathbf{x} = [5, 14, 6]$?

$$\begin{aligned} P(y=+, \mathbf{x}=[5, 14, 6]) &= P(y=+) P([5, 14, 6] | y=+) \\ &= \frac{3}{7} \times P(\mathbf{x}=[5, 14, 6] | \mathbf{x} \sim \mathcal{N}(\hat{\mu}_+, \hat{\Sigma})) \\ &= \frac{3}{7} \times \frac{1}{(2\pi)^{3/2} |\hat{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \hat{\mu}_+)^T \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\mu}_+) \right] \\ &= 16.63E-11 \end{aligned}$$

$$\begin{aligned} P(y=-, \mathbf{x}=[5, 14, 6]) &= P(y=-) P([5, 14, 6] | y=-) \\ &= \frac{4}{7} \times \frac{1}{(2\pi)^{3/2} |\hat{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \hat{\mu}_-)^T \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\mu}_-) \right] \\ &= 43.33E-11 \end{aligned}$$

$$\begin{aligned} \bullet \quad P(y=+ | [5, 14, 6]) &= \\ \frac{P(y=+, [5, 14, 6])}{P(y=+, [5, 14, 6]) + P(y=-, [5, 14, 6])} &= 0.2774 \\ P(y=- | [5, 14, 6]) &= 0.7226 \end{aligned}$$

LDA learns an LTU

- Consider 2-class case with a 0/1 loss function
- Classify $\hat{y} = 1$ if

$$\log \frac{P(y=1|x)}{P(y=0|x)} > 0 \quad \text{iff} \quad \log \frac{P(y=1, x)}{P(y=0, x)} > 0$$

$$\begin{aligned}\frac{P(x, y=1)}{P(x, y=0)} &= \frac{P(y=1) \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1)\right]}{P(y=0) \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x - \mu_0)^\top \Sigma^{-1}(x - \mu_0)\right]} \\ &= \frac{P(y=1) \exp\left[-\frac{1}{2}(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1)\right]}{P(y=0) \exp\left[-\frac{1}{2}(x - \mu_0)^\top \Sigma^{-1}(x - \mu_0)\right]}\end{aligned}$$

$$\ln \frac{P(x, y=1)}{P(x, y=0)} = \ln \frac{P(y=1)}{P(y=0)} - \frac{1}{2} [(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1) - (x - \mu_0)^\top \Sigma^{-1}(x - \mu_0)]$$

LDA Learns an LTU (2)

- $$\begin{aligned} & (\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_0) \\ &= \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} + \\ &\quad \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 \end{aligned}$$
- As $\boldsymbol{\Sigma}^{-1}$ is symmetric,
$$\begin{aligned} \dots &= 2 \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 \\ \Rightarrow \ln \frac{P(\mathbf{x}, y=1)}{P(\mathbf{x}, y=0)} &= \\ \ln \frac{P(y=1)}{P(y=0)} - \frac{1}{2} & [(\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_0)] \\ = \ln \frac{P(y=1)}{P(y=0)} &+ \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \frac{1}{2} \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 - \frac{1}{2} \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \\ = \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) &+ \ln \frac{P(y=1)}{P(y=0)} + \frac{1}{2} \boldsymbol{\mu}_0^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 - \frac{1}{2} \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \end{aligned}$$

LDA Learns an LTU (3)

$$\ln \frac{P(\mathbf{x}, y=1)}{P(\mathbf{x}, y=0)} = \mathbf{x}^\top \Sigma^{-1} (\mu_1 - \mu_0) + \ln \frac{P(y=1)}{P(y=0)} + \frac{1}{2} \mu_0^\top \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^\top \Sigma^{-1} \mu_1$$

- So let...

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_0)$$

$$c = \ln \frac{P(y=1)}{P(y=0)} + \frac{1}{2} \mu_0^\top \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^\top \Sigma^{-1} \mu_1$$

- Classify $\hat{y} = 1$ iff $\mathbf{w} \cdot \mathbf{x} + c > 0$

LTU!!

Variants of LDA

- Covariance matrix Σ
- n features; k classes

Same for all classes?	Diagonal	#param's	Name
+	+	k	
+	-	n^2	LDA
-	+	$k n$	Naïve Gaussian Classifier
-	-	$k n^2$	General Gaussian Classifier

Generalizations of LDA

■ General Gaussian Classifier

Allow each class k to have its own Σ_k

⇒ Classifier = quadratic threshold unit (not LTU)

■ Naïve Gaussian Classifier

Allow each class k to have its own Σ_k

but require each Σ_k be diagonal.

⇒ within each class,

any pair of features x_i and x_j are independent

□ Classifier is still quadratic threshold unit
but with a restricted form

Summary of Linear Discriminant Analysis

- **Learns Joint Probability Distr'n $P(y, x)$**
- **Direct Computation.**
MLEstimate of $P(y, x)$ computed directly from data without search.
But need to invert matrix, which is $O(n^3)$
- **Eager:**
Classifier constructed from training examples, which can then be discarded.
- **Batch:** Only a batch algorithm.
An online LDA alg requires online alg for incrementally updating Σ^{-1}
[Easy if Σ^{-1} is diagonal. . .]

Two Geometric Views of LDA

View 1: Mahalanobis Distance

- Squared Mahalanobis distance between \mathbf{x} and $\boldsymbol{\mu}$

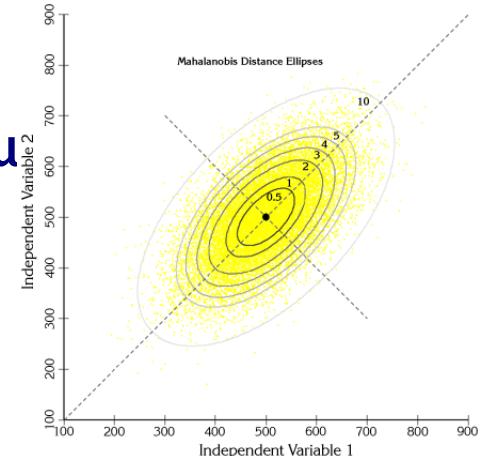
$$D_M^2(\mathbf{x}, \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$\boldsymbol{\Sigma}^{-1} \approx$ linear distortion
... converts standard Euclidean distance into Mahalanobis distance.

- LDA classifies \mathbf{x} as 0 if

$$D_M^2(\mathbf{x}, \boldsymbol{\mu}_0) < D_M^2(\mathbf{x}, \boldsymbol{\mu}_1)$$

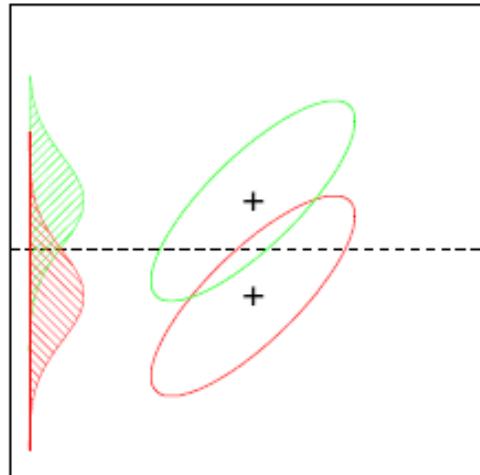
- $\log P(\mathbf{x} | y = k) \approx \log \pi_k - \frac{1}{2} D_M^2(\mathbf{x}, \boldsymbol{\mu}_k)$



View 2: Most Informative Low Dimensional Projection

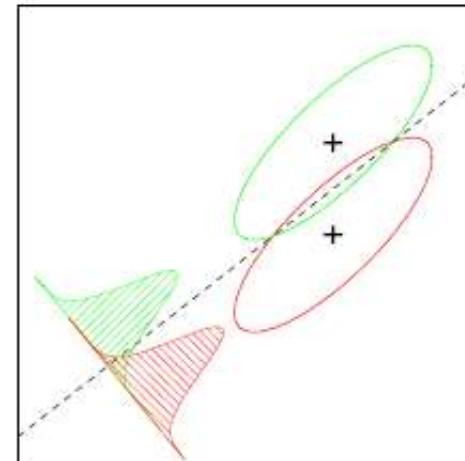
- LDA
 - Finds $K-1$ dim hyperplane ($K = \text{number of classes}$)
 - Project \mathbf{x} and $\{\mu_k\}$ to that hyperplane
 - Classify \mathbf{x} as nearest μ_k within hyperplane
- Goal: Hyperplane that maximally separates projection of \mathbf{x} 's wrt Σ^{-1}

\approx Fisher's Linear Discriminant



project onto:

vertical axis



LDA's w

Fisher Linear Discriminant

- Recall any vector \mathbf{w} projects $\Re^n \rightarrow \Re$
- Goal: Want \mathbf{w} that “separates” classes
 - Each $\mathbf{w} \cdot \mathbf{x}^+$ far from each $\mathbf{w} \cdot \mathbf{x}^-$

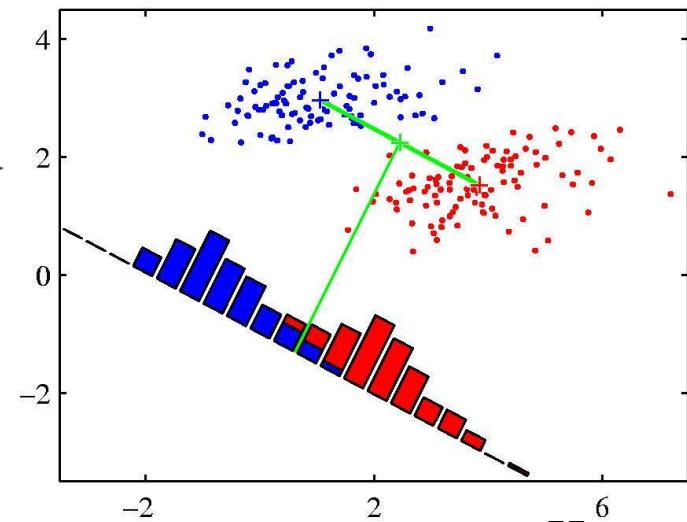
- Using $\mathbf{m}_+ = \frac{\sum_i y^{(i)} \cdot \mathbf{x}^{(i)}}{\sum_i y^{(i)}}$ $\mathbf{m}_- = \frac{\sum_i (1-y^{(i)}) \cdot \mathbf{x}^{(i)}}{\sum_i (1-y^{(i)})}$

Mean of \mathbf{x} 's projections:

$$m_+ = \frac{\sum_i y^{(i)} \mathbf{w}^\top \cdot \mathbf{x}^{(i)}}{\sum_i y^{(i)}} = \mathbf{w}^\top \cdot \mathbf{m}_+$$

$$m_- = \frac{\sum_i (1-y^{(i)}) \mathbf{w}^\top \cdot \mathbf{x}^{(i)}}{\sum_i (1-y^{(i)})} = \mathbf{w}^\top \cdot \mathbf{m}_-$$

- Perhaps project on $\mathbf{m}_+ - \mathbf{m}_-$?
- Still overlap... why?

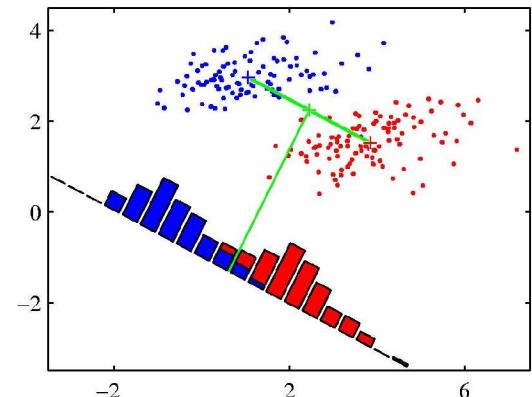


Fisher Linear Discriminant

- Using $m_+ = \frac{\sum_i y^{(i)} \cdot x^{(i)}}{\sum_i y^{(i)}}$ $m_- = \frac{\sum_i (1-y^{(i)}) \cdot x^{(i)}}{\sum_i (1-y^{(i)})}$

Mean of x's projections:

$$m_+ = \frac{\sum_i y^{(i)} w^\top \cdot x^{(i)}}{\sum_i y^{(i)}} = w^\top \cdot m_+$$
$$m_- = \frac{\sum_i (1-y^{(i)}) w^\top \cdot x^{(i)}}{\sum_i (1-y^{(i)})} = w^\top \cdot m_-$$



- Problem with $m_+ - m_-$:
- Doesn't consider "scatter" within class
- Goal: Want w that "separates" classes
 - Each $w \cdot x^+$ far from each $w \cdot x^-$
 - Positive x^+ 's: $w \cdot x^+$ close to each other
 - Negative x^- 's: $w \cdot x^-$ close to each other
- "scatter" of +instance; -instance
 - $s_+^2 = \sum_i y^{(i)} (w \cdot x^{(i)} - m_+)^2$
 - $s_-^2 = \sum_i (1 - y^{(i)}) (w \cdot x^{(i)} - m_+)^2$

Fisher Linear Discriminant

- Recall any vector \mathbf{w} projects $\Re^n \rightarrow \Re$
- Goal: Want \mathbf{w} that “separates” classes
 - Positive \mathbf{x}^+ 's: $\mathbf{w} \cdot \mathbf{x}^+$ close to each other
 - Negative \mathbf{x}^- 's: $\mathbf{w} \cdot \mathbf{x}^-$ close to each other
 - Each $\mathbf{w} \cdot \mathbf{x}^+$ far from each $\mathbf{w} \cdot \mathbf{x}^-$

- Using $\mathbf{m}_+ = \frac{\sum_i y^{(i)} \cdot \mathbf{x}^{(i)}}{\sum_i y^{(i)}}$ $\mathbf{m}_- = \frac{\sum_i (1-y^{(i)}) \cdot \mathbf{x}^{(i)}}{\sum_i (1-y^{(i)})}$

Mean of \mathbf{x} 's projections:

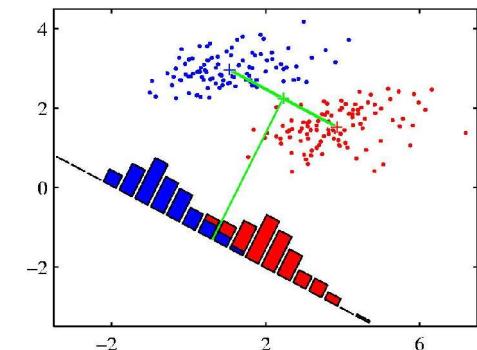
$$m_+ = \frac{\sum_i y^{(i)} \mathbf{w}^\top \cdot \mathbf{x}^{(i)}}{\sum_i y^{(i)}} = \mathbf{w}^\top \cdot \mathbf{m}_+$$

$$m_- = \frac{\sum_i (1-y^{(i)}) \mathbf{w}^\top \cdot \mathbf{x}^{(i)}}{\sum_i (1-y^{(i)})} = \mathbf{w}^\top \cdot \mathbf{m}_-$$

- “scatter” of $+$ instance; $-$ instance

- $s_+^2 = \sum_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} - m_+)^2$

- $s_-^2 = \sum_i (1 - y^{(i)}) (\mathbf{w} \cdot \mathbf{x}^{(i)} - m_+)^2$



FLD, con't

- Separate means \mathbf{m}_- and \mathbf{m}_+

$$\Rightarrow \text{maximize } (\mathbf{m}_+ - \mathbf{m}_-)^2$$

- Minimize each spread $\mathbf{s}_+^2, \mathbf{s}_-^2$

$$\Rightarrow \text{maximize } (\mathbf{s}_+^2 + \mathbf{s}_-^2)$$

- Objective function: maximize

$$J_S(w) = \frac{(\mathbf{m}_+ - \mathbf{m}_-)^2}{(\mathbf{s}_+^2 + \mathbf{s}_-^2)}$$

$$\begin{aligned} \#1: (\mathbf{m}_+ - \mathbf{m}_-)^2 &= (\mathbf{w}^\top \mathbf{m}_+ - \mathbf{w}^\top \mathbf{m}_-)^2 \\ &= \mathbf{w}^\top (\mathbf{m}_+ - \mathbf{m}_-) (\mathbf{m}_+ - \mathbf{m}_-)^T \mathbf{w} = \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \end{aligned}$$

“between-class scatter”

$$\mathbf{S}_B = (\mathbf{m}_+ - \mathbf{m}_-) (\mathbf{m}_+ - \mathbf{m}_-)^T$$

$$J_S(w) = \frac{(m_+ - m_-)^2}{(s_+^2 + s_-^2)}$$

FLD, III

- $\mathbf{s}_+^2 = \sum_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} - m_+)^2$
 $= \sum_i \mathbf{w}^\top y^{(i)} (\mathbf{x}^{(i)} - m_+) (\mathbf{x}^{(i)} - m_+)^\top \mathbf{w}$
 $= \mathbf{w}^\top \mathbf{S}_+ \mathbf{w}$

$$\mathbf{S}_+ = \sum_i y^{(i)} (\mathbf{x}^{(i)} - m_+) (\mathbf{x}^{(i)} - m_+)^\top$$

... “within-class scatter matrix” for +

$$\mathbf{S}_- = \sum_i (1 - y^{(i)}) (\mathbf{x}^{(i)} - m_-) (\mathbf{x}^{(i)} - m_-)^\top$$

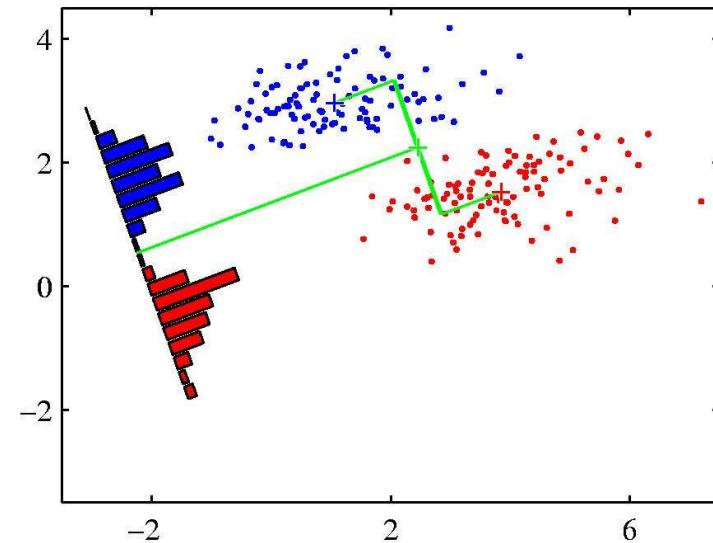
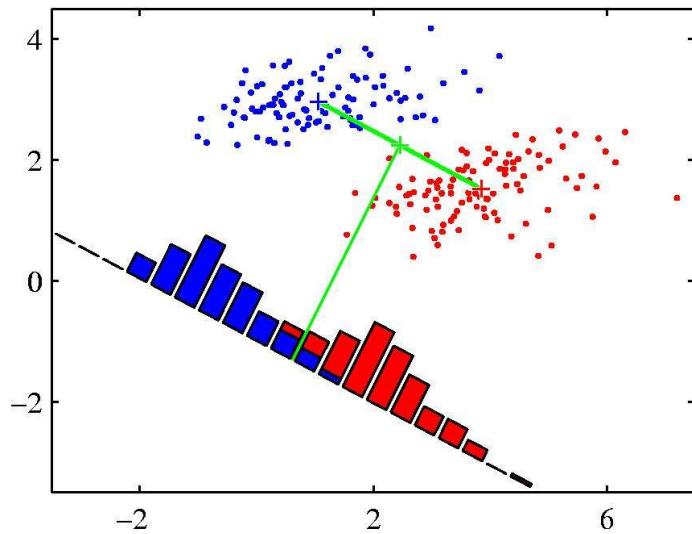
... “within-class scatter matrix” for -

- $\mathbf{S}_w = \mathbf{S}_+ + \mathbf{S}_-$ so $\mathbf{s}_+^2 + \mathbf{s}_-^2 = \mathbf{w}^\top \mathbf{S}_w \mathbf{w}$

FLD, IV

$$J_S(w) = \frac{(m_+ - m_-)^2}{(s_+^2 + s_-^2)} = \frac{w^T S_B w}{w^T S_w w} = \frac{(w^T(m_1 - m_2))^2}{w^T S_w w}$$

- Solving $\frac{\partial J_S(w)}{\partial w_j} = 0 \Rightarrow \mathbf{w} = \alpha \mathbf{S}_B^{-1}(\mathbf{m}_+ - \mathbf{m}_-)$



FLD, V

$$J_S(w) = \frac{(m_+ - m_-)^2}{(s_+^2 + s_-^2)} = \frac{w^T S_B w}{w^T S_w w} = \frac{(w^T(m_1 - m_2))^2}{w^T S_w w}$$

- Solving $\frac{\partial J_S(w)}{\partial w_j} = 0 \Rightarrow w = \alpha S_B^{-1}(m_+ - m_-)$
- When $P(x | y_i) \sim N(\mu_i; \Sigma)$
 \exists LINEAR DISCRIMINANT: $w = \Sigma^{-1}(\mu_+ - \mu_-)$
 \Rightarrow FLD is optimal classifier,
 if classes normally distributed
- Can use even if not normal:
 After projecting d -dim to 1,
 just use any classification method
- Analogous derivation for $K > 2$ classes

Comparing LMS, Logistic Regression, LDA

- Which is best: *LMS*, *LR*, *LDA* ?
- Ongoing debate within machine learning community about relative merits of
 - direct classifiers [*LMS*]
 - conditional models $P(y | \mathbf{x})$ [*LR*]
 - generative models $P(y, \mathbf{x})$ [*LDA*]
- Stay tuned...



Issues in Debate

- **Statistical efficiency**

If generative model $P(y, x)$ is correct, then ...
usually gives better accuracy, particularly if training sample is small

- **Computational efficiency**

Generative models typically easiest to compute
(LDA computed directly, without iteration)

- **Robustness to changing loss functions**

LMS must re-train the classifier when the loss function changes.
... no retraining for generative and conditional models

- **Robustness to model assumptions.**

Generative model usually performs poorly when the assumptions are violated.

Eg, LDA works poorly if $P(x | y)$ is non-Gaussian.

Logistic Regression is more robust, ... LMS is even more robust

- **Robustness to missing values and noise.**

In many applications, some of the features x_{ij} may be missing or corrupted for some of the training examples.

Generative models typically provide better ways of handling this than non-generative models.

Other Algorithms for learning LTUs

- **Naive Bayes** [Discuss later]

For $K = 2$ classes, produces LTU

- **Winnow** [?Discuss later?]

Can handle large numbers of “irrelevant” features

- (features whose weights should be zero)