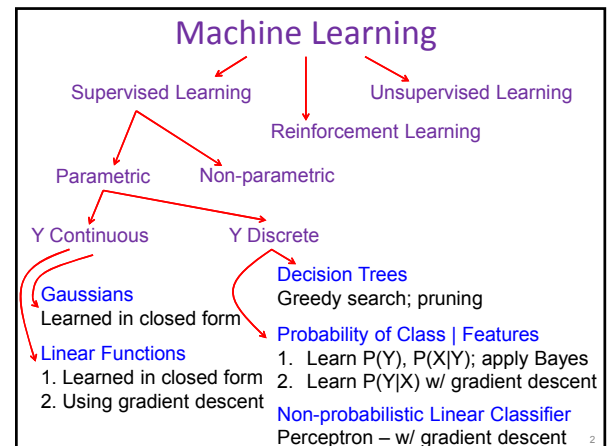


# CSE 446 Perceptron Learning Winter 2012

Dan Weld

Some slides from Carlos Guestrin, Luke Zettlemoyer



## Hypothesis Expressiveness

### Linear

- Naïve Bayes
- Logistic Regression
- **Perceptron**
- Support Vector Machines

### Nonlinear

- **Decision Trees**
- Neural Networks
- **Ensembles**
- **Kernel Methods**
- **Nearest Neighbor**
- Graphical Models

3

## Logistic Regression

- **Want to Learn:**  $h: X \mapsto Y$ 
  - $X$  – features
  - $Y$  – target classes
- **Probabilistic Discriminative Classifier**
  - Assume some **functional form** for  $P(Y|X)$ 
    - Logistic Function
    - Accepts both discrete & continuous features
  - Estimate parameters of  $P(Y|X)$  directly from training data
  - This is the **'discriminative' model**
    - Directly learn  $P(Y|X)$
    - But **cannot generate a sample of the data**,
    - No way to compute  $P(X)$

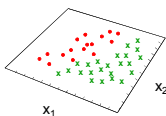
4

## Earthquake or Nuclear Test?

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$



linear  
classification  
rule!

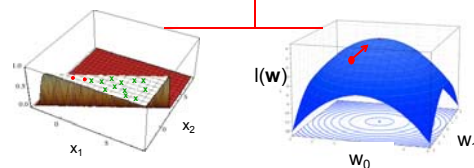
5

## Logistic w/ Initial Weights

$$w_0 = 20 \quad w_1 = -5 \quad w_2 = 10$$

$$\text{Loss}(H_w) = \text{Error}(H_w, \text{data})$$

$$\text{Minimize Error} \rightarrow \text{Maximize } l(w) = \ln P(D_Y | D_X, H_w)$$



Update rule:

$$\Delta w = \eta \nabla_w l(w)$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(w)}{\partial w_i}$$

6

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w})] \right\}$$

- Assume functional form for
  - $P(Y|X)$  **no assumptions**
  - Est params from training data
- Handles discrete & cont features
- Directly calculate  $P(Y|X=x)$ 
  - Can't generate data sample

**Cool!!!!**

- Yes, but only in Gaussian case

Linear function!  
Coefficients  
expressed with  
original Gaussian  
parameters!

### Derive form for $P(Y|X)$ for continuous $X_i$

$$P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$$

$$= \frac{1}{1 + \exp\left(\ln \frac{1-\theta}{\theta} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

$$w_0 = \ln \frac{1-\theta}{\theta} + \frac{\mu_{i0}^2 + \mu_{i1}^2}{2\sigma_i^2}$$

$$w_i = \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2}$$

### Gaussian Naïve Bayes vs. Logistic Regression

Set of Gaussian Naïve Bayes parameters (feature variance independent of class label)



Set of Logistic Regression parameters

Can go both ways, we only did one way

- **Representation equivalence**
  - But only in a special case!!! (GNB with class-independent variances)
- But what's the difference???
- **LR makes no assumptions about  $P(X|Y)$  in learning!!!**
- **Loss function!!!**
  - Optimize different functions ! Obtain different solutions

### Naïve Bayes vs. Logistic Regression

Consider  $Y$  boolean,  $X = \langle X_1 \dots X_n \rangle$  continuous

Number of parameters:

- Naïve Bayes:  $4n+1$
- Logistic Regression:  $n+1$

$$P(Y) = \theta$$

$$P(X|Y) = N(\mu, \sigma)$$

$$P(X|\neg Y) = N(\mu', \sigma')$$

Estimation method:

- Naïve Bayes parameter estimates are uncoupled
- Logistic Regression parameter estimates are coupled

### Naïve Bayes vs. Logistic Regression

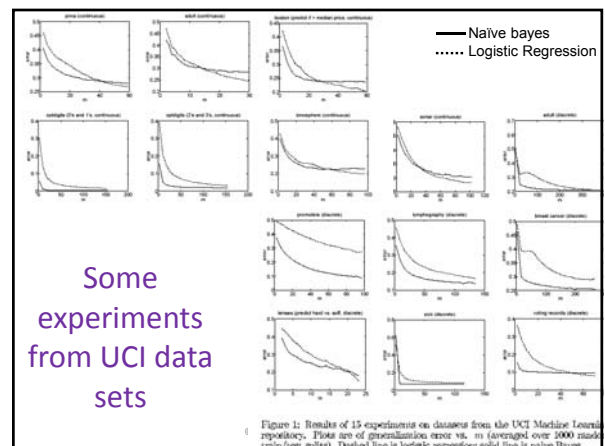
[Ng & Jordan, 2002]

- Generative vs. Discriminative classifiers
- Asymptotic comparison (# training examples  $\rightarrow$  infinity)
  - when model correct
    - GNB (with class independent variances) and LR produce identical classifiers
  - when model incorrect
    - LR is less biased – does not assume conditional independence
      - therefore LR expected to outperform GNB

### Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

- Generative vs. Discriminative classifiers
- Non-asymptotic analysis
  - convergence rate of parameter estimates, ( $n = \#$  of attributes in  $X$ )
    - Size of training data to get close to infinite data solution
    - Naïve Bayes needs  $O(\log n)$  samples
    - Logistic Regression needs  $O(n)$  samples
  - GNB converges more quickly to its (perhaps less helpful) asymptotic estimates



## What you should know about Logistic Regression (LR)

- Gaussian Naïve Bayes with class-independent variances representationally equivalent to LR
  - Solution differs because of objective (loss) function
- In general, NB and LR make different assumptions
  - NB: Features independent given class ! assumption on  $P(\mathbf{X}|Y)$
  - LR: Functional form of  $P(Y|\mathbf{X})$ , no assumption on  $P(\mathbf{X}|Y)$
- LR is a linear classifier
  - decision rule is a hyperplane
- LR optimized by conditional likelihood
  - no closed-form solution
  - concave ! global optimum with gradient ascent
  - Maximum conditional a posteriori corresponds to regularization
- Convergence rates
  - GNB (usually) needs less data
  - LR (usually) gets to better solutions in the limit

## Perceptrons

20

## Who needs probabilities?

- Previously: model data with distributions
- Joint:  $P(\mathbf{X}, Y)$ 
  - e.g. Naïve Bayes
- Conditional:  $P(Y|\mathbf{X})$ 
  - e.g. Logistic Regression
- But wait, why probabilities?
- Lets try to be error-driven!

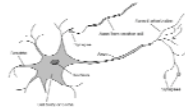
| mpg  | cylinders | displacement | horsepower | weight | acceleration | modelyear | make |
|------|-----------|--------------|------------|--------|--------------|-----------|------|
| good | 4         | 97           | 75         | 2265   | 18.2         | 77        | asia |
| bad  | 6         | 199          | 90         | 2648   | 15           | 70        | amer |
| bad  | 4         | 121          | 110        | 2000   | 12.8         | 77        | euro |
| bad  | 8         | 350          | 175        | 4100   | 13           | 73        | amer |
| bad  | 6         | 198          | 95         | 3102   | 16.5         | 74        | amer |
| bad  | 4         | 106          | 94         | 2379   | 16.5         | 73        | asia |
| bad  | 4         | 113          | 95         | 2228   | 14           | 71        | asia |
| bad  | 8         | 302          | 139        | 3570   | 12.8         | 78        | amer |
| ...  | ...       | ...          | ...        | ...    | ...          | ...       | ...  |
| good | 4         | 120          | 79         | 2625   | 18.6         | 80        | amer |
| bad  | 8         | 455          | 225        | 4425   | 10           | 70        | amer |
| good | 4         | 107          | 86         | 2464   | 15.5         | 76        | euro |
| bad  | 5         | 131          | 103        | 2930   | 15.9         | 78        | euro |

## Generative vs. Discriminative

- Generative classifiers:
  - E.g. naïve Bayes
  - A joint probability model with evidence variables
  - Query model for causes given evidence
- Discriminative classifiers:
  - No generative model, no Bayes rule, often no probabilities at all!
  - Try to predict the label  $Y$  directly from  $X$
  - Robust, accurate with varied features
  - Loosely: **mistake driven rather than model driven**

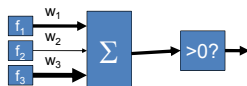
## Linear Classifiers

- Inputs are feature values
- Each feature has a weight
- Sum is the activation



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
  - Positive, output class 1
  - Negative, output class 2



## Example: Spam

- Imagine 3 features (spam is “positive” class):

- free (number of occurrences of “free”)
- money (occurrences of “money”)
- BIAS (intercept, always has value 1)

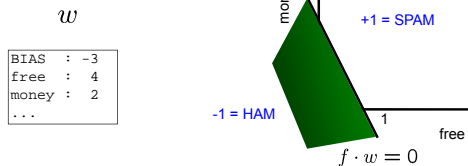
$$w \cdot f(x) = \sum_i w_i \cdot f_i(x)$$

| $x$       | $f(x)$    | $w$ |             |
|-----------|-----------|-----|-------------|
| BIAS : 1  | BIAS : -3 |     | $(1)(-3) +$ |
| free : 1  | free : 4  |     | $(1)(4) +$  |
| money : 1 | money : 2 |     | $(1)(2) +$  |
| ...       | ...       |     | ...         |
|           |           |     | $= 3$       |

$w \cdot f(x) > 0 \rightarrow \text{SPAM!!!}$

## Binary Decision Rule

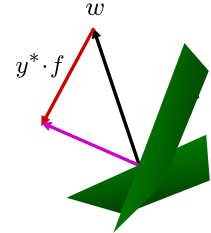
- In the space of feature vectors
  - Examples are points
  - Any weight vector is a hyperplane
  - One side corresponds to  $Y=+1$
  - Other corresponds to  $Y=-1$



## Binary Perceptron Algorithm

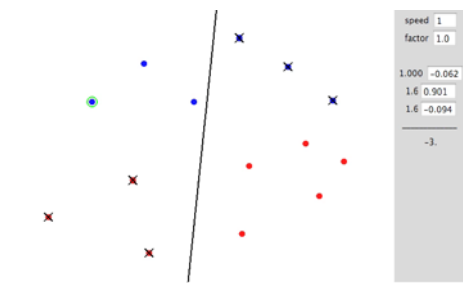
- Start with zero weights
  - For each training instance  $(x, y^*)$ :
    - Classify with current weights
- $$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$
- If correct (i.e.,  $y=y^*$ ), no change!
  - If wrong: update

$$w = w + y^* \cdot f$$



## Examples: Perceptron

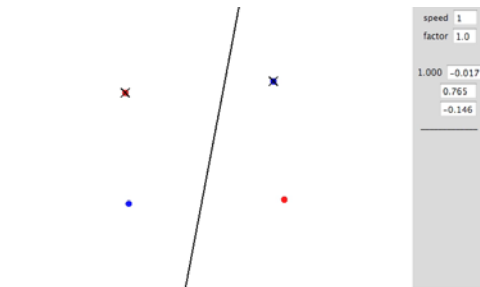
- Separable Case



[http://isl.ira.uka.de/neuralNetCourse/2004/VL\\_11\\_5/Perceptron.html](http://isl.ira.uka.de/neuralNetCourse/2004/VL_11_5/Perceptron.html)

## Examples: Perceptron

- Inseparable Case



[http://isl.ira.uka.de/neuralNetCourse/2004/VL\\_11\\_5/Perceptron.html](http://isl.ira.uka.de/neuralNetCourse/2004/VL_11_5/Perceptron.html)

## From Logistic Regression to the Perceptron: 2 easy steps!

- Logistic Regression: (in vector notation):  $y$  is  $\{0,1\}$

$$w = w + \eta \sum_j [y_j^* - p(y_j^* | x_j, w)] f(x_j)$$

- Perceptron:  $y$  is  $\{0,1\}$ ,  $y(x;w)$  is prediction given  $w$

$$w = w + [y^* - y(x;w)] f(x)$$

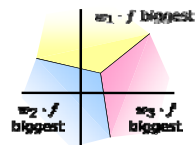
### Differences?

- Drop the  $\sum_j$  over training examples: **online vs. batch learning**
- Drop the dist'n: **probabilistic vs. error driven learning**

## Multiclass Decision Rule

- If we have more than two classes:

- Have a weight vector for each class:  $w_y$
- Calculate an activation for each class



$$\text{activation}_w(x, y) = w_y \cdot f(x)$$

- Highest activation wins

$$y = \arg \max_y (\text{activation}_w(x, y))$$

## Example

“win the vote”  
 “win the election”  
 “win the game”

$w_{SPORTS}$

|      |   |  |
|------|---|--|
| BIAS | : |  |
| win  | : |  |
| game | : |  |
| vote | : |  |
| the  | : |  |
| ...  | : |  |

$w_{POLITICS}$

|      |   |  |
|------|---|--|
| BIAS | : |  |
| win  | : |  |
| game | : |  |
| vote | : |  |
| the  | : |  |
| ...  | : |  |

$w_{TECH}$

|      |   |  |
|------|---|--|
| BIAS | : |  |
| win  | : |  |
| game | : |  |
| vote | : |  |
| the  | : |  |
| ...  | : |  |

## Example

“win the vote”



|      |   |   |
|------|---|---|
| BIAS | : | 1 |
| win  | : | 1 |
| game | : | 0 |
| vote | : | 1 |
| the  | : | 1 |
| ...  | : |   |

$w_{SPORTS}$

|      |   |    |
|------|---|----|
| BIAS | : | -2 |
| win  | : | 4  |
| game | : | 4  |
| vote | : | 0  |
| the  | : | 0  |
| ...  | : |    |

$w_{POLITICS}$

|      |   |   |
|------|---|---|
| BIAS | : | 1 |
| win  | : | 2 |
| game | : | 0 |
| vote | : | 4 |
| the  | : | 0 |
| ...  | : |   |

$w_{TECH}$

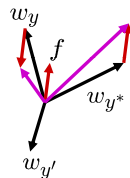
|      |   |   |
|------|---|---|
| BIAS | : | 2 |
| win  | : | 0 |
| game | : | 2 |
| vote | : | 0 |
| the  | : | 0 |
| ...  | : |   |

## The Multi-class Perceptron Alg.

- Start with zero weights
- Iterate training examples
  - Classify with current weights
 
$$y = \arg \max_y w_y \cdot f(x)$$

$$= \arg \max_y \sum_i w_{y,i} \cdot f_i(x)$$
  - If correct, no change!
  - If wrong: lower score of wrong answer, raise score of right answer
 
$$w_y = w_y - f(x)$$

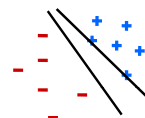
$$w_{y^*} = w_{y^*} + f(x)$$



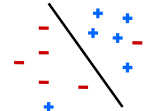
## Properties of Perceptrons

- Separability: some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the margin or degree of separability

Separable



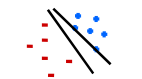
Non-Separable



$$\text{mistakes} < \frac{k}{\delta^2}$$

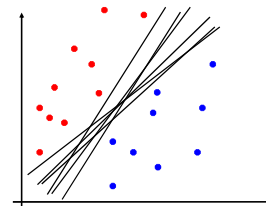
## Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)
- Mediocre generalization: finds a “barely” separating solution
- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining is a kind of overfitting



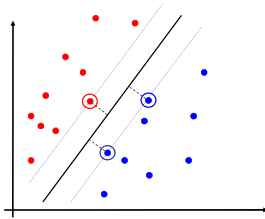
## Linear Separators

- Which of these linear separators is optimal?



## Support Vector Machines

- Maximizing the margin: good according to intuition, theory, practice
- SVMs find the separator with max margin

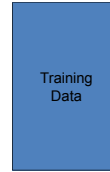


SVM

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w_y \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

## Three Views of Classification



- Naïve Bayes:
  - Parameters from data statistics
  - Parameters: probabilistic interpretation
  - Training: one pass through the data
- Logistic Regression:
  - Parameters from gradient ascent
  - Parameters: linear, probabilistic model, and discriminative
  - Training: one pass through the data per gradient step; regularization essential
- The Perceptron:
  - Parameters from reactions to mistakes
  - Parameters: discriminative interpretation
  - Training: go through the data until held-out accuracy maxes out