

Neural Networks

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 14th, 2007

©2005-2007 Carlos Guestrin

1

Logistic regression

- $P(Y|X)$ represented by:

$$P(Y = 1 | x, W) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}} = g(w_0 + \sum_i w_i x_i)$$

- Learning rule – MLE:

Compute derivative

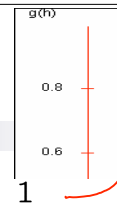
$$\begin{aligned} \frac{\partial \ell(W)}{\partial w_i} &= \sum_j x_i^j [y^j - P(Y^j = 1 | x^j, W)] \\ &= \sum_j x_i^j [y^j - g(w_0 + \sum_i w_i x_i^j)] \end{aligned}$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j \delta^j$$

learning rate

$$\delta^j = y^j - g(w_0 + \sum_i w_i x_i^j)$$

truth prediction



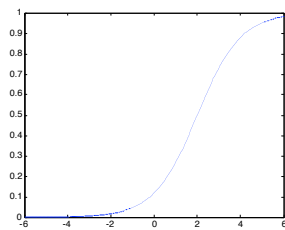
©2005-2007 Carlos Guestrin

2

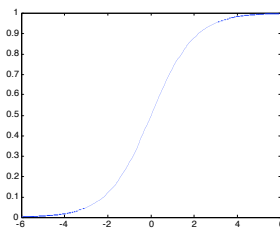
Sigmoid

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

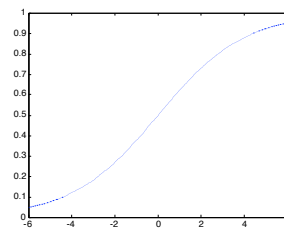
$w_0=2, w_1=1$



$w_0=0, w_1=1$



$w_0=0, w_1=0.5$

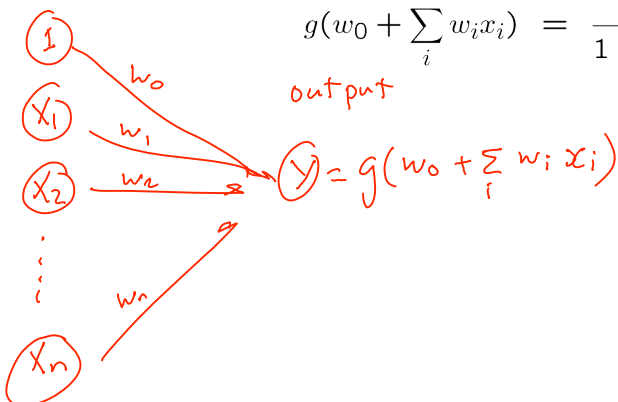


©2005-2007 Carlos Guestrin

3

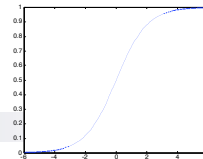
Perceptron as a graph

input node



$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

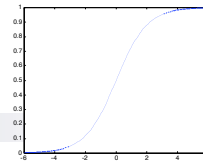
output



©2005-2007 Carlos Guestrin

4

Linear perceptron classification region



(+)
 $w_0 + \sum_i w_i x_i \geq 0$

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

$w_0 + \sum_i w_i x_i < 0$

(-)

$w_0 + \sum_i w_i x_i < 0$

©2005-2007 Carlos Guestrin

5

Optimizing the perceptron

- Trained to minimize sum-squared error

$$\ell(W) = \frac{1}{2} \sum_{j=1}^m [y^j - g(w_0 + \sum_i w_i x_i^j)]^2$$

Annotations:
 - y^j : datapoints
 - y^j : truth
 - $g(w_0 + \sum_i w_i x_i^j)$: prediction
 - 2 : squared

©2005-2007 Carlos Guestrin

6

Derivative of sigmoid

$$\frac{\partial \ell(W)}{\partial w_i} = - \sum_j [y^j - g(w_0 + \sum_i w_i x_i^j)] x_i^j g'(w_0 + \sum_i w_i x_i^j)$$

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

©2005-2007 Carlos Guestrin

7

The perceptron learning rule

$$w_i \leftarrow w_i + \eta \sum_j x_i^j \delta^j$$

$$\delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)] g'(w_0 + \sum_i w_i x_i^j)$$

$$g^j = g(w_0 + \sum_i w_i x_i^j)$$

■ Compare to MLE:

$$w_i \leftarrow w_i + \eta \sum_j x_i^j \delta^j \quad \delta^j = [y^j - g(w_0 + \sum_i w_i x_i^j)]$$

©2005-2007 Carlos Guestrin

8

Perceptron, linear classification, Boolean functions

- Can learn $x_1 \vee x_2$
- Can learn $x_1 \wedge x_2$
- Can learn any conjunction or disjunction

©2005-2007 Carlos Guestrin

9

Perceptron, linear classification, Boolean functions

- Can learn majority
- Can perceptrons do everything?

©2005-2007 Carlos Guestrin

10

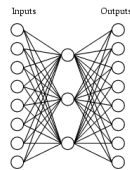
Going beyond linear classification

- Solving the XOR problem

Hidden layer

- Perceptron: $out(\mathbf{x}) = g(w_0 + \sum_i w_i x_i)$
- 1-hidden layer:
$$out(\mathbf{x}) = g\left(w_0 + \sum_k w_k g(w_0^k + \sum_i w_i^k x_i)\right)$$

Example data for NN with hidden layer



A target function:

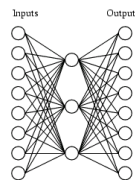
Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

Can this be learned??

13

Learned weights for hidden layer

A network:

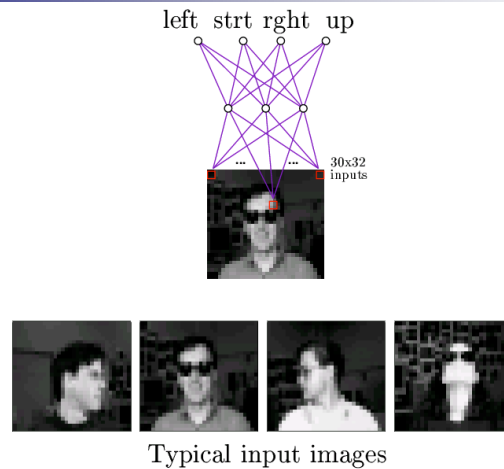


Learned hidden layer representation:

Input	Hidden Values	Output
10000000	→ .89 .04 .08 →	10000000
01000000	→ .01 .11 .88 →	01000000
00100000	→ .01 .97 .27 →	00100000
00010000	→ .99 .97 .71 →	00010000
00001000	→ .03 .05 .02 →	00001000
00000100	→ .22 .99 .99 →	00000100
00000010	→ .80 .01 .98 →	00000010
00000001	→ .60 .94 .01 →	00000001

14

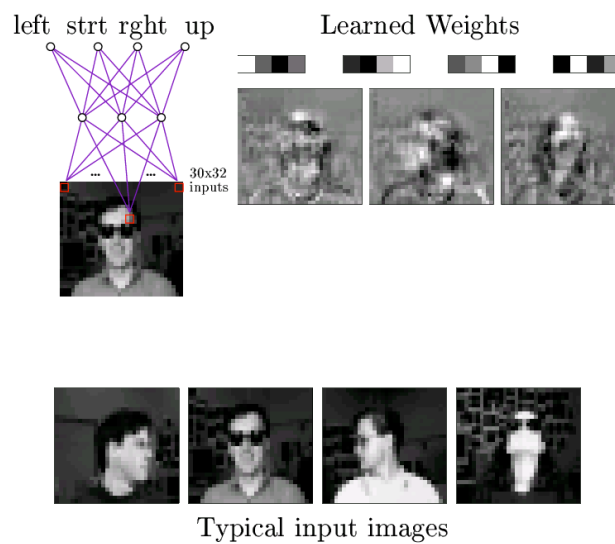
NN for images



90% accurate learning head pose, and recognizing 1-of-20 faces

15

Weights in NN for images



©2005-2007 Carlos Guestrin

16

Forward propagation for 1-hidden layer - Prediction

- 1-hidden layer:

$$out(\mathbf{x}) = g \left(w_0 + \sum_k w_k g(w_0^k + \sum_i w_i^k x_i) \right)$$

©2005-2007 Carlos Guestrin

17

Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_k}$

$$\ell(W) = \frac{1}{2} \sum_j [y^j - out(\mathbf{x}^j)]^2$$

Dropped w_0 to make derivation simpler

$$out(\mathbf{x}) = g \left(\sum_{k'} w_{k'} g(\sum_{i'} w_{i'}^{k'} x_{i'}) \right)$$

$$\frac{\partial \ell(W)}{\partial w_k} = \sum_{j=1}^m [y - out(\mathbf{x}^j)] \frac{\partial out(\mathbf{x}^j)}{\partial w_k}$$

©2005-2007 Carlos Guestrin

18

Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_i^k}$

$$\ell(W) = \frac{1}{2} \sum_j [y^j - \text{out}(\mathbf{x}^j)]^2$$

Dropped w_0 to make derivation simpler

$$\text{out}(\mathbf{x}) = g \left(\sum_{k'} w_{k'} g \left(\sum_{i'} w_{i'}^{k'} x_{i'} \right) \right)$$

$$\frac{\partial \ell(W)}{\partial w_i^k} = \sum_{j=1}^m -[y - \text{out}(\mathbf{x}^j)] \frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_i^k}$$

Multilayer neural networks

Forward propagation – prediction

- Recursive algorithm
- Start from input layer
- Output of node V_k with parents U_1, U_2, \dots :

$$V_k = g\left(\sum_i w_i^k U_i\right)$$

©2005-2007 Carlos Guestrin

21

Back-propagation – learning

- Just gradient descent!!!
- Recursive algorithm for computing gradient
- For each example
 - Perform forward propagation
 - Start from output layer
 - Compute gradient of node V_k with parents U_1, U_2, \dots
 - Update weight w_i^k

©2005-2007 Carlos Guestrin

22

Many possible response functions

- Sigmoid
- Linear
- Exponential
- Gaussian
- ...

©2005-2007 Carlos Guestrin

23

Convergence of backprop

- Perceptron leads to convex optimization
 - Gradient descent reaches **global minima**
- Multilayer neural nets **not convex**
 - Gradient descent gets stuck in local minima
 - Hard to set learning rate
 - Selecting number of hidden units and layers = fuzzy process
 - NNs falling in disfavor in last few years
 - We'll see later in semester, *kernel trick* is a good alternative
 - Nonetheless, neural nets are one of the most used ML approaches

©2005-2007 Carlos Guestrin

24

Training set error

- Neural nets represent complex functions
 - Output becomes more complex with gradient steps

©2005-2007 Carlos Guestrin

25

Overfitting

- Output fits training data “too well”
 - Poor test set accuracy
- Overfitting the training data
 - Related to bias-variance tradeoff
 - One of central problems of ML
- Avoiding overfitting?
 - More training data
 - Regularization
 - Early stopping

©2005-2007 Carlos Guestrin

26

What you need to know about neural networks

- Perceptron:
 - ☐ Representation
 - ☐ Perceptron learning rule
 - ☐ Derivation
- Multilayer neural nets
 - ☐ Representation
 - ☐ Derivation of backprop
 - ☐ Learning rule
- Overfitting
 - ☐ Definition
 - ☐ Training set versus test set
 - ☐ Learning curve

©2005-2007 Carlos Guestrin

27

Instance-based Learning

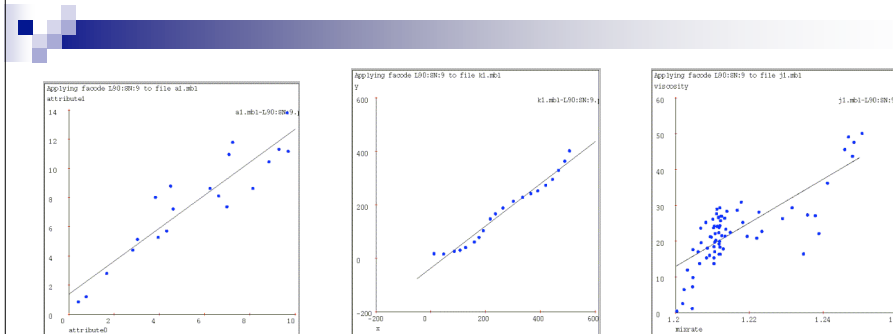
Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

February 14th, 2007

©2005-2007 Carlos Guestrin

28

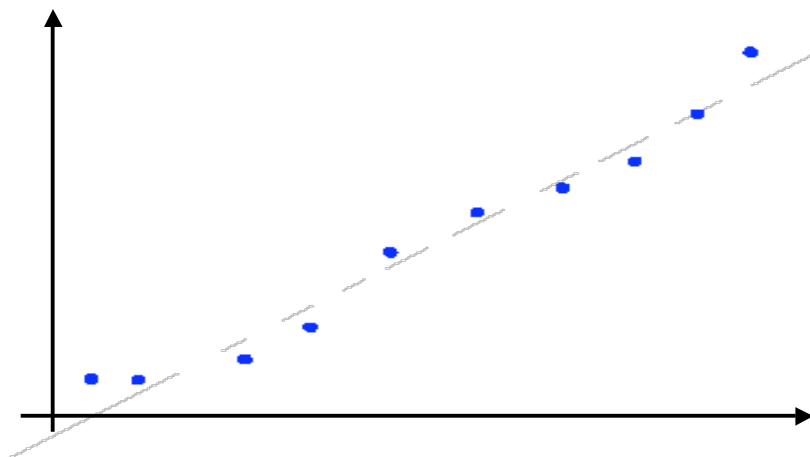
Why not just use Linear Regression?



©2005-2007 Carlos Guestrin

29

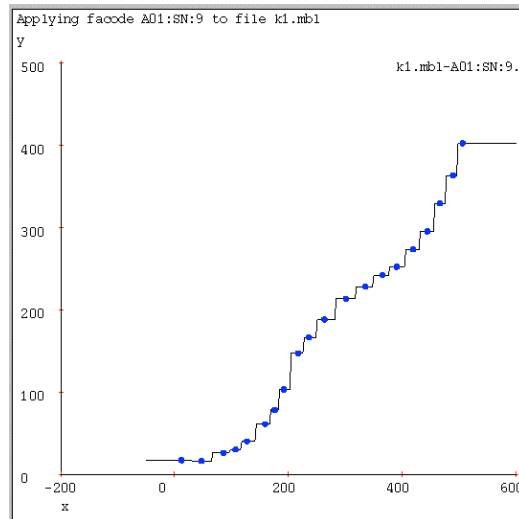
Using data to predict new data



©2005-2007 Carlos Guestrin

30

Nearest neighbor



©2005-2007 Carlos Guestrin

31

Univariate 1-Nearest Neighbor

Given datapoints $(x_1, y_1) (x_2, y_2) \dots (x_N, y_N)$, where we assume $y_i = f(x_i)$ for some unknown function f .

Given query point x_q , your job is to predict $\hat{y} \approx f(x_q)$

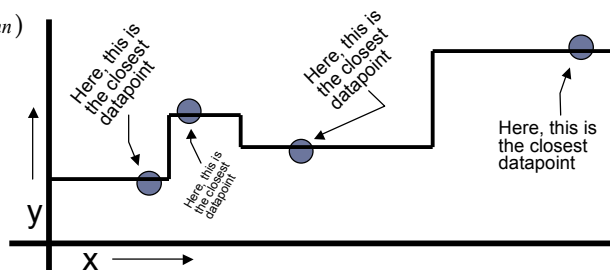
Nearest Neighbor:

1. Find the closest x_i in our set of datapoints

$$i(nn) = \underset{i}{\operatorname{argmin}} |x_i - x_q|$$

2. Predict $\hat{y} = y_{i(nn)}$

Here's a dataset with one input, one output and four datapoints.



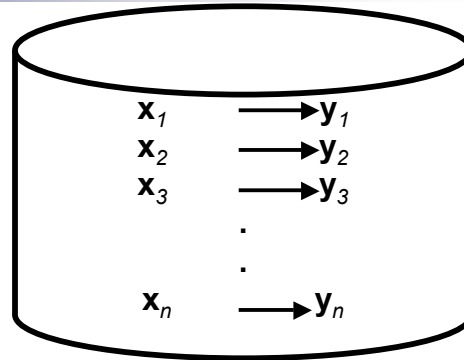
©2005-2007 Carlos Guestrin

32

1-Nearest Neighbor is an example of.... Instance-based learning

A function approximator that has been around since about 1910.

To make a prediction, search database for similar datapoints, and fit with the local points.



Four things make a memory based learner:

- A distance metric
- How many nearby neighbors to look at?
- A weighting function (optional)
- How to fit with the local points?

©2005-2007 Carlos Guestrin

33

1-Nearest Neighbor

Four things make a memory based learner:

1. *A distance metric*
Euclidian (and many more)
2. *How many nearby neighbors to look at?*
One
3. *A weighting function (optional)*
Unused
4. *How to fit with the local points?*
Just predict the same output as the nearest neighbor.

©2005-2007 Carlos Guestrin

34

Multivariate 1-NN examples

Regression

Classification

©2005-2007 Carlos Guestrin

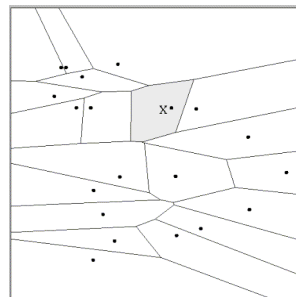
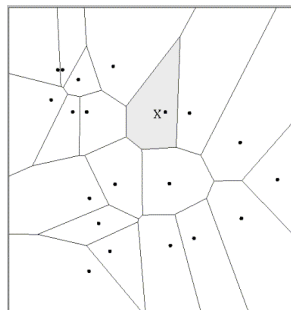
35

Multivariate distance metrics

Suppose the input vectors x_1, x_2, \dots, x_n are two dimensional:

$\mathbf{x}_1 = (x_{11}, x_{12}), \mathbf{x}_2 = (x_{21}, x_{22}), \dots, \mathbf{x}_N = (x_{N1}, x_{N2})$.

One can draw the nearest-neighbor regions in input space.



$$Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 \quad Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (3x_{i2} - 3x_{j2})^2$$

The relative scalings in the distance metric affect region shapes.

©2005-2007 Carlos Guestrin

36