



Simple Model Selection Cross Validation Regularization Neural Networks

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 13th, 2007

©2005-2007 Carlos Guestrin

1

OK... now we'll learn to pick those
darned parameters...

- **Selecting features (or basis functions)**

- ☐ Linear regression
- ☐ Naïve Bayes
- ☐ Logistic regression

- **Selecting parameter value**

- ☐ Prior strength
 - Naïve Bayes, linear and logistic regression
- ☐ Regularization strength
 - Naïve Bayes, linear and logistic regression
- ☐ Decision trees
 - MaxpChance, depth, number of leaves
- ☐ Boosting
 - Number of rounds

- More generally, these are called **Model Selection** Problems

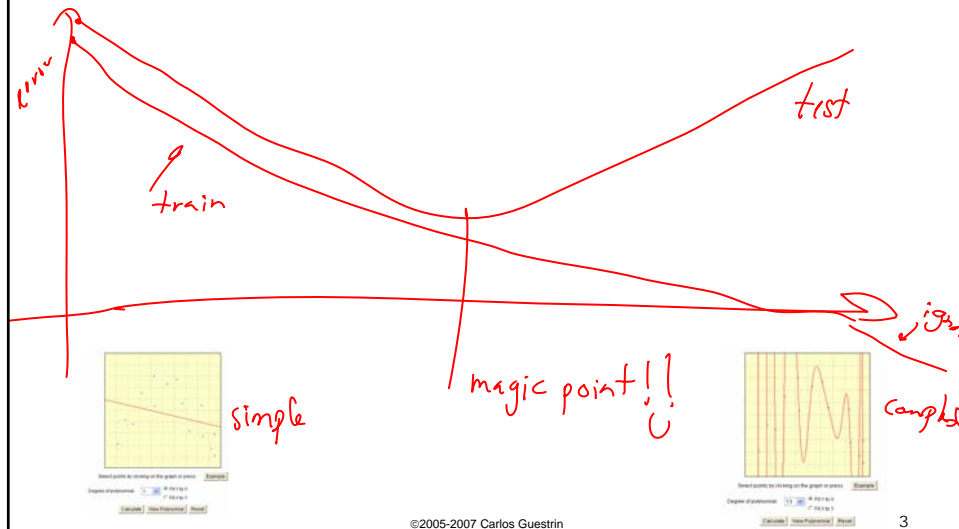
- Today:

- ☐ Describe basic idea
- ☐ Introduce very important concept for tuning learning approaches: **Cross-Validation**

©2005-2007 Carlos Guestrin

2

Test set error as a function of model complexity



3

Simple greedy model selection algorithm

- Pick a dictionary of features $\{1, x, x^2, x^3, x^4, \dots\}$
 - e.g., polynomials for linear regression
- Greedy heuristic:
 - Start from empty (or simple) set of features $F_0 = \emptyset$ eg., $F_0 = \{1, x\}$
 - Run learning algorithm for current set of features F_t
 - Obtain h_t
 - Select **next best feature** X_j
 - e.g., X_j that results in lowest training error learner when learning with $F_t \cup \{X_j\}$
 - $F_{t+1} \leftarrow F_t \cup \{X_j\}$
 - Recurse

©2005-2007 Carlos Guestrin

4

Greedy model selection

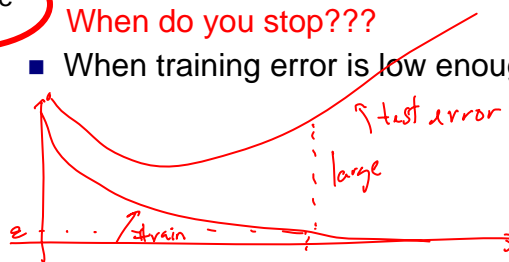
- Applicable in many settings:
 - Linear regression: Selecting basis functions
 - Naïve Bayes: Selecting (independent) features $P(X_i|Y)$
 - Logistic regression: Selecting features (basis functions)
 - Decision trees: Selecting leaves to expand
- Only a heuristic!
 - But, sometimes you can prove something cool about it
 - e.g., [Krause & Guestrin '05]: Near-optimal in some settings that include Naïve Bayes
- There are many more elaborate methods out there

©2005-2007 Carlos Guestrin

5

Simple greedy model selection algorithm

- Greedy heuristic:
 - ...
 - Select **next best feature** X_j
 - e.g., X_j that results in lowest training error learner when learning with $F_t \cup \{X_j\}$
 - $F_{t+1} \leftarrow F_t \cup \{X_j\}$
 - Recurse
- When training error is low enough? $\leq \epsilon$



©2005-2007 Carlos Guestrin

6

Simple greedy model selection algorithm

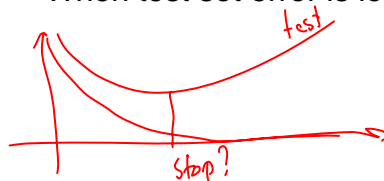
Greedy heuristic:

- ...
- Select **next best feature** X_i
 - e.g., X_i that results in lowest training error learner when learning with $F_t \cup \{X_i\}$
- $F_{t+1} \leftarrow F_t \cup \{X_i\}$
- **Recurse**

When do you stop???

- ~~When training error is low enough?~~
- When test set error is low enough?

never, never!!
;



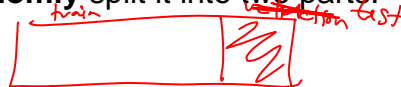
©2005-2007 Carlos Guestrin

7

Validation set

Thus far: Given a dataset, **randomly** split it into two parts:

- Training data – $\{x_1, \dots, x_{N_{\text{train}}}\}$
- Test data – $\{x_1, \dots, x_{N_{\text{test}}}\}$

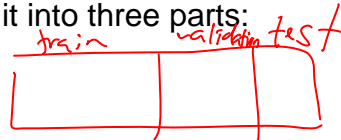


But **Test data must always remain independent!**

- Never ever ever ever learn on test data, including for model selection

Given a dataset, **randomly** split it into three parts:

- Training data – $\{x_1, \dots, x_{N_{\text{train}}}\}$
- Validation data – $\{x_1, \dots, x_{N_{\text{valid}}}\}$
- Test data – $\{x_1, \dots, x_{N_{\text{test}}}\}$



Use validation data for tuning learning algorithm, e.g., model selection

- Save test data for very final evaluation

©2005-2007 Carlos Guestrin

8

Simple greedy model selection algorithm

- Greedy heuristic:
 - ...
 - Select **next best feature** X_i
 - e.g., X_i that results in lowest training error learner when learning with $F_t \cup \{X_i\}$
 - $E_{t+1} < E_t \cup \{X_i\}$
 - Recurse
- train
val.
test
- | | | |
|--|--|--|
| | | |
|--|--|--|
- When do you stop???
 - ~~When training error is low enough?~~
 - ~~When test set error is low enough?~~
 - When validation set error is low enough?

©2005-2007 Carlos Guestrin

9

Simple greedy model selection algorithm

- Greedy heuristic:
 - ...
 - Select **next best feature** X_i
 - e.g., X_i that results in lowest training error learner when learning with $F_t \cup \{X_i\}$
 - $E_{t+1} < E_t \cup \{X_i\}$
 - Recurse
- over fit
to validation sets
- When do you stop???
 - ~~When training error is low enough?~~
 - ~~When test set error is low enough?~~
 - ~~When validation set error is low enough?~~
 - Man!!! OK, should I just repeat until I get tired???
 - I am tired now...
 - No, "There is a better way!"

©2005-2007 Carlos Guestrin

10

(LOO) Leave-one-out cross validation

- Consider a **validation set with 1 example**:

- D – training data
- D_i – training data with i th data point moved to validation set

- Learn classifier h_{D_i} with D_i dataset

- Estimate true error as:

- 0 if h_{D_i} classifies i th data point correctly
- 1 if h_{D_i} is wrong about i th data point
- Seems really bad estimator, but wait!

$E[\mathbb{1}(h_{D_i}(x_i) \neq y_i)] = \text{error}_{\text{true}}$
 \rightarrow unbiased estimate of true error of h_{D_i}

- LOO cross validation**: Average over all data points i :

- For each data point you leave out, learn a new classifier h_{D_i}
- Estimate error as:

$$\text{error}_{LOO} = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(h_{D \setminus i}(x^i) \neq y^i)$$

\nwarrow train examples

does h_{D_i} predict x_i correctly

©2005-2007 Carlos Guestrin

11

LOO cross validation is (almost) unbiased estimate of true error!

- When computing **LOOCV error**, we only use $m-1$ data points

- So it's not estimate of true error of learning with m data points!
- Usually pessimistic, though – learning with less data typically gives worse answer

- LOO is almost unbiased!**

- Let $\text{error}_{\text{true}, m-1}$ be true error of learner when you only get $m-1$ data points
- In homework, you'll prove that LOO is unbiased estimate of $\text{error}_{\text{true}, m-1}$:

$$E_D[\text{error}_{LOO}] = \text{error}_{\text{true}, m-1}$$

- Great news!**

- Use LOO error for model selection!!!

©2005-2007 Carlos Guestrin

12

Simple greedy model selection algorithm

Greedy heuristic:

- ...
 - Select **next best feature** X_i
 - e.g., X_i that results in lowest training error learner when learning with $F_t \cup \{X_i\}$
 - $F_{t+1} \leftarrow F_t \cup \{X_i\}$
 - Recurse
- When do you stop???
- ~~When training error is low enough?~~
 - ~~When test set error is low enough?~~
 - ~~When validation set error is low enough?~~
 - **STOP WHEN error_{LOO} IS LOW!!!**

©2005-2007 Carlos Guestrin

13

Using LOO error for model selection



©2005-2007 Carlos Guestrin

14

Computational cost of LOO

- Suppose you have 100,000 data points
- You implemented a great version of your learning algorithm
 - Learns in only 1 second
- Computing LOO will take about 1 day!!!
 - If you have to do for each choice of basis functions, it will take fooooooreeeve'!!!
- Solution 1: Preferred, but not usually possible
 - Find a cool trick to compute LOO (e.g., see homework)

©2005-2007 Carlos Guestrin

15

Solution 2 to complexity of computing LOO:

(More typical) **Use k -fold cross validation**

- Randomly **divide training data into k equal parts**
 - D_1, \dots, D_k
- For each i
 - Learn classifier $h_{D \setminus D_i}$ using data point not in D_i m/k points
 - Estimate error of $h_{D \setminus D_i}$ on validation set D_i :

$$error_{D_i} = \frac{1}{m} \sum_{(x^j, y^j) \in D_i} \mathbb{I}(h_{D \setminus D_i}(x^j) \neq y^j)$$

size of data \rightarrow m \leftarrow *error on left out section D_i*
- **k -fold cross validation error is average over data splits:**

$$error_{k\text{-fold}} = \frac{1}{k} \sum_{i=1}^k error_{D_i}$$

- k -fold cross validation properties:
 - **Much faster to compute** than LOO
 - **More (pessimistically) biased** – using much less data, only $m(k-1)/k$
 - **Usually, $k = 10$** ☺

©2005-2007 Carlos Guestrin

16

Regularization – Revisited

■ Model selection 1: **Greedy**

- Pick subset of features that have yield low LOO error

■ Model selection 2: **Regularization**

- Include all possible features! $\{1, x, \dots, x^m\}$
- Penalize “complicated” hypothesis

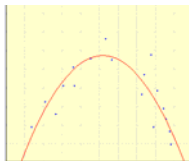
©2005-2007 Carlos Guestrin

17

Regularization in linear regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$



$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



- Regularized least-squares (a.k.a. ridge regression), for $\lambda \geq 0$:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{\sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2}_{\text{training error}} + \underbrace{\lambda \sum_{i=1}^k w_i^2}_{\substack{\text{penalty for large } w_i \\ \text{regularization term}}}$$

©2005-2007 Carlos Guestrin

18

Other regularization examples

Logistic regression regularization

- Maximize data likelihood minus **penalty for large parameters**

$$\arg \max_{\mathbf{w}} \underbrace{\sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w})}_{\text{conditional log likelihood}} - \underbrace{\lambda \sum_i w_i^2}_{\text{regularization}}$$

- Biases towards small parameter values**

Naïve Bayes regularization (smoothing)

- Prior** over likelihood of features
- Biases away from zero probability** outcomes

Decision tree regularization

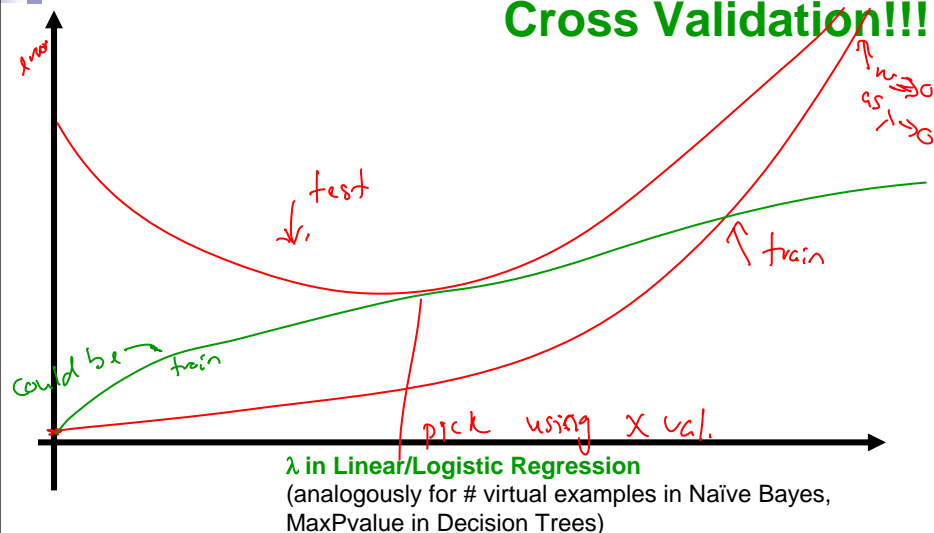
- Many possibilities, e.g., **Chi-Square test** and **MaxPvalue** parameter
- Biases towards smaller trees**

©2005-2007 Carlos Guestrin

19

How do we pick magic parameter?

Cross Validation!!!!



©2005-2007 Carlos Guestrin

20

Regularization and Bayesian learning

- $p(\mathbf{w} | Y, \mathbf{X}) \propto P(Y | \mathbf{X}, \mathbf{w})p(\mathbf{w})$
- $\ln p(\mathbf{w} | D) = \text{constant} + \ln P(Y | \mathbf{X}, \mathbf{w}) + \ln p(\mathbf{w})$
- We already saw that **regularization for logistic regression** corresponds to **MAP for zero mean, Gaussian prior for \mathbf{w}**

$p(\mathbf{w}) \propto \exp\left(-\frac{\mathbf{w}^2}{\sigma^2}\right)$

regularization = placing a prior over parameters
 - Similar interpretation for other learning approaches:
 - **Linear regression:** Also zero mean, Gaussian prior for \mathbf{w}
 - **Naïve Bayes:** Directly defined as prior over parameters (Smoothing)
 - **Decision trees:** Trickier to define... but we'll get back to this

©2005-2007 Carlos Guestrin

21

Occam's Razor



- William of Ockham (1285-1349) Principle of Parsimony:
 - "One should not increase, beyond what is necessary, the number of entities required to explain anything."
- Regularization penalizes for "complex explanations"
- Alternatively (but pretty much the same), use Minimum Description Length (MDL) Principle:
 - minimize $\text{length}(\text{misclassifications}) + \text{length}(\text{hypothesis})$

length = number of bits you got wrong

length = how hard to write down

size of tree

complexity model

minimize sum
- $\text{length}(\text{misclassifications})$ – e.g., #wrong training examples
- $\text{length}(\text{hypothesis})$ – e.g., size of decision tree

©2005-2007 Carlos Guestrin

22

Minimum Description Length Principle

- MDL prefers small hypothesis that fit data well:

$$h_{MDL} = \arg \min_h L_{C_1}(\mathcal{D} | h) + L_{C_2}(h)$$

- $L_{C_1}(\mathcal{D} | h)$ – description length of data under code C_1 given h
 - Only need to describe points that h doesn't explain (classify correctly)
- $L_{C_2}(h)$ – description length of hypothesis h

- Decision tree example

- $L_{C_1}(\mathcal{D} | h)$ – #bits required to describe data given h
 - If all points correctly classified, $L_{C_1}(\mathcal{D} | h) = 0$
- $L_{C_2}(h)$ – #bits necessary to encode tree
- Trade off quality of classification with tree size

©2005-2007 Carlos Guestrin

23

Bayesian interpretation of MDL Principle

MAP estimate

$$h_{MAP} = \arg \max_h [P(\mathcal{D} | h) P(h)]$$

take log

$$= \arg \max_h [\log_2 P(\mathcal{D} | h) + \log_2 P(h)]$$

argmax = argmin - f

$$= \arg \min_h [-\log_2 P(\mathcal{D} | h) - \log_2 P(h)]$$

likelihood prior

- **Information theory fact:**

- Smallest code for event of probability p requires $-\log_2 p$ bits

- MDL interpretation of MAP:

- $-\log_2 P(\mathcal{D} | h)$ – length of \mathcal{D} under hypothesis h think about # bit required to represent wrong egs
- $-\log_2 P(h)$ – length of hypothesis h (there is hidden parameter here)
- MAP prefers simpler hypothesis:
 - minimize $length(misclassifications) + length(hypothesis)$ MAP is also MDL

- **In general, Bayesian approach usually looks for simpler hypothesis** – Acts as a regularizer

©2005-2007 Carlos Guestrin

24

What you need to know about Model Selection, Regularization and Cross Validation

- Cross validation
 - (Mostly) Unbiased estimate of true error
 - LOOCV is great, but hard to compute
 - k -fold much more practical
 - Use for selecting parameter values!
- Model selection
 - Search for a model with low cross validation error
- Regularization
 - Penalizes for complex models
 - Select parameter with cross validation
 - Really a Bayesian approach
- Minimum description length
 - Information theoretic interpretation of regularization
 - Relationship to MAP

©2005-2007 Carlos Guestrin

25

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

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

learning rate

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

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

truth prediction

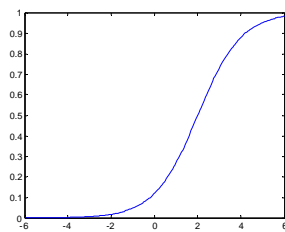
©2005-2007 Carlos Guestrin

26

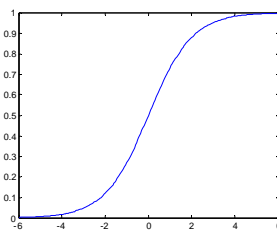
Sigmoid

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

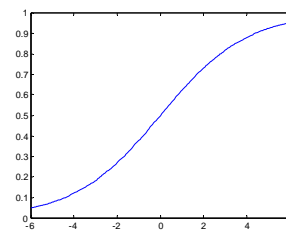
$w_0=2, w_1=1$



$w_0=0, w_1=1$



$w_0=0, w_1=0.5$

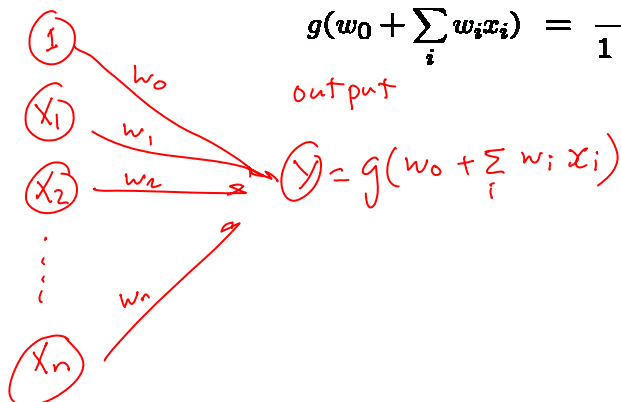


©2005-2007 Carlos Guestrin

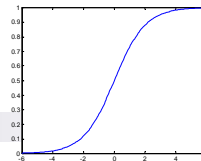
27

Perceptron as a graph

input node



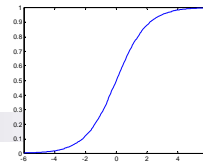
$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$



©2005-2007 Carlos Guestrin

28

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

29

Optimizing the perceptron

- Trained to minimize sum-squared error

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

data points → y^j (truth) $g(w_0 + \sum_i w_i x_i^j)$ (prediction)

squared

©2005-2007 Carlos Guestrin

30

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

31

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

32

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

33

Perceptron, linear classification, Boolean functions

- Can learn majority
- Can perceptrons do everything?

©2005-2007 Carlos Guestrin

34

Going beyond linear classification

- Solving the XOR problem

©2005-2007 Carlos Guestrin

35

Hidden layer

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

©2005-2007 Carlos Guestrin

36

Example data for NN with hidden layer



A target function:

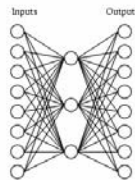
Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

Can this be learned??

37

Learned weights for hidden layer

A network:

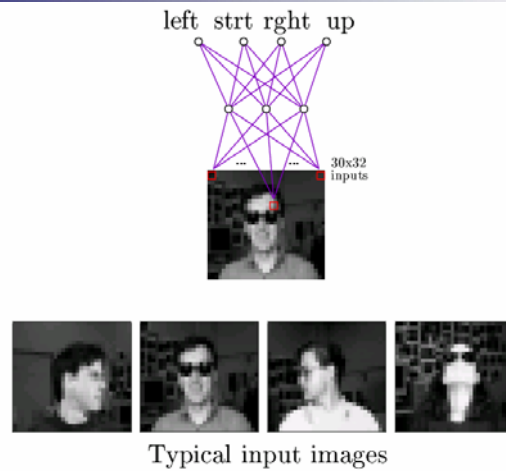


Learned hidden layer representation:

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

38

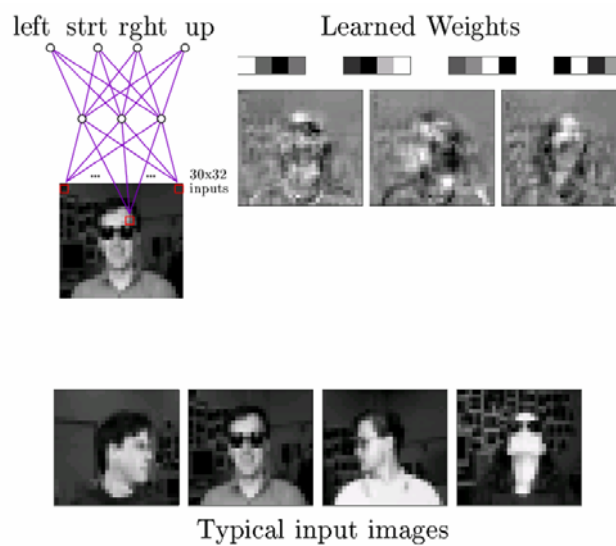
NN for images



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

39

Weights in NN for images



©2005-2007 Carlos Guestrin

40