

Data Structures for Fast Kernel Regression and Locally Weighted Regression

The Auton Lab
Carnegie Mellon University



www.autonlab.org

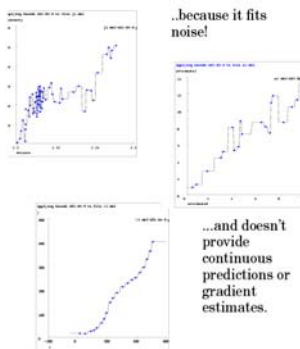
What is LWPR?

Locally weighted polynomial regression is an algorithm for learning non-linear numeric functions.

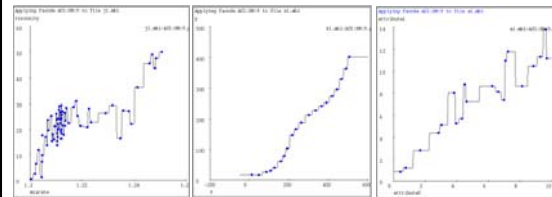
Here, we will...

- Review LWPR (travelling via Nearest Neighbor and Kernel Regression).
- Praise its virtues.
- Mourn its computational expense.
- Discuss recently proposed methods for making it faster.
- Introduce and evaluate a new multiresolution data structure.

Why not just use 1-nearest neighbor?

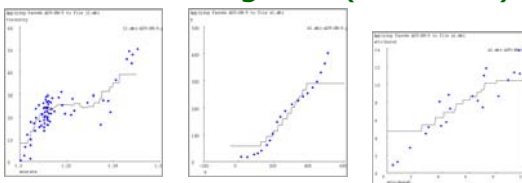


Why not just use one-Nearest Neighbor?



- Discontinuities
- Fits the noise

k-Nearest Neighbor (here k=9)



A magnificent job of noise-smoothing. Three cheers for 9-nearest-neighbor. But the lack of gradients and the jerkiness isn't good.

Appalling behavior! Loses all the detail that join-the-dots and 1-nearest-neighbor gave us, yet smears the ends.

Fits much less of the noise, captures trends. But still, frankly, pathetic compared with linear regression.

K-nearest neighbor for function fitting smooths away noise, but there are clear deficiencies.

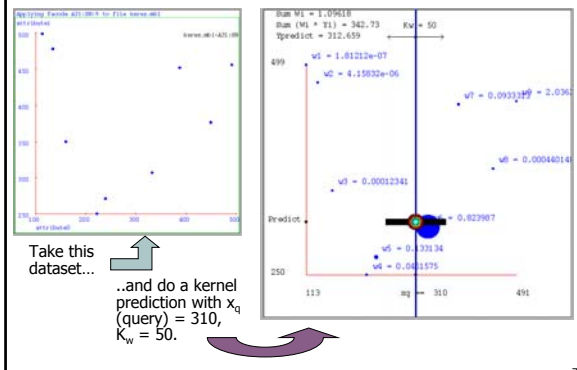
What can we do about all the discontinuities that k-NN gives us?

Kernel Regression

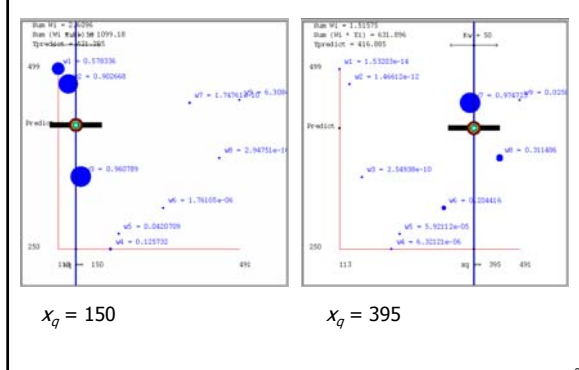
Four things make a memory based learner:

1. A distance metric
Scaled Euclidian
2. How many nearby neighbors to look at?
All of them
3. A weighting function (optional)
 $w_i = \exp(-D(x_{\text{query}})^2 / K_w^2)$
Nearby points to the query are weighted strongly, far points weakly. The K_w parameter is the **Kernel Width**. Very important.
4. How to fit with the local points?
Predict the weighted average of the outputs:
 $\text{predict} = \sum w_i y_i / \sum w_i$

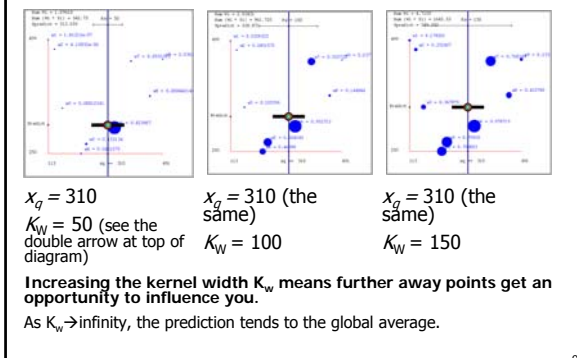
Kernel Regression in Pictures



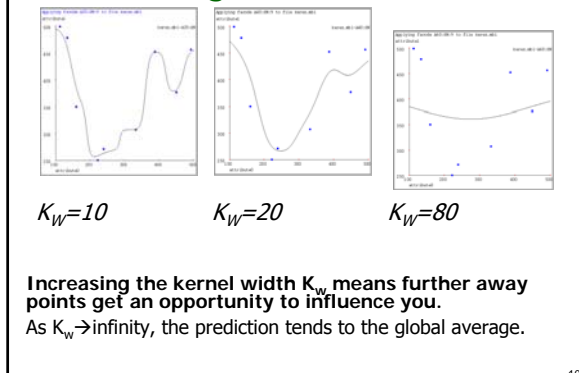
Varying the Query



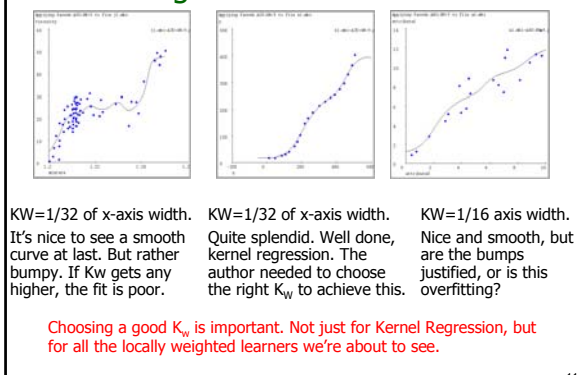
Varying the kernel width



Kernel Regression Predictions



Kernel Regression on our test cases



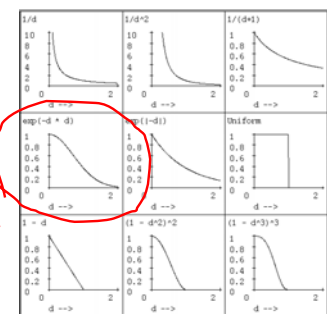
Weighting functions

Let

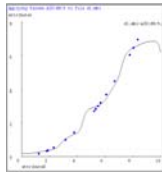
$$d = D(x_i, x_{\text{query}}) / K_w$$

Then here are some commonly used weighting functions...

(we use a Gaussian)



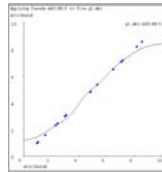
Kernel Regression can look bad



KW = Best.

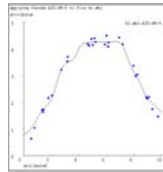
Clearly not capturing the simple structure of the data.. Note the complete failure to extrapolate at edges.

Time to try something more powerful...



KW = Best.

Also much too local. Why wouldn't increasing Kw help? Because then it would all be "smeared".



KW = Best.

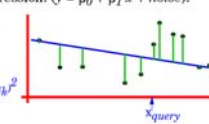
Three noisy linear segments. But best kernel regression gives poor gradients.

Locally weighted linear regression

Global Regression: ($y = \beta_0 + \beta_1 x + \text{noise}$):

Minimizes sum of squared residuals:

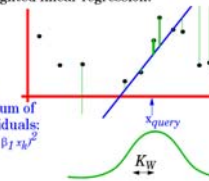
$$\sum_k (y_k - \beta_0 - \beta_1 x_k)^2$$



Locally weighted linear regression:

Minimizes weighted sum of squared residuals:

$$\sum_k w_k (y_k - \beta_0 - \beta_1 x_k)^2$$



Locally Weighted Regression

Four things make a memory-based learner:

1. A distance metric

Scaled Euclidian

2. How many nearby neighbors to look at?

All of them

3. A weighting function (optional)

$$w_k = \exp(-D(x_{iw}, x_{query})^2 / K_w^2)$$

Nearby points to the query are weighted strongly, far points weakly. The K_w parameter is the Kernel Width.

4. How to fit with the local points?

First form a local linear model. Find the β that minimizes the locally weighted sum of squared residuals:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{k=1}^N w_k^2 (y_k - \beta^T x_k)^2 \quad \text{Then predict } y_{\text{predict}} = \beta^T x_{\text{query}}$$

How LWR works

1. For each point (x_k, y_k) compute w_k

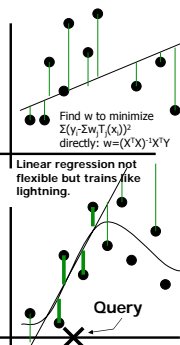
2. Let $WX = \operatorname{Diag}(w_1, \dots, w_N)X$

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix} \rightarrow \begin{bmatrix} w_1 & w_1 x_{11} & w_1 x_{12} & \dots & w_1 x_{1D} \\ w_2 & w_2 x_{21} & w_2 x_{22} & \dots & w_2 x_{2D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_N & w_N x_{N1} & w_N x_{N2} & \dots & w_N x_{ND} \end{bmatrix}$$

$$X \rightarrow WX$$

3. Let $WY = \operatorname{Diag}(w_1, \dots, w_N)Y$, so that $y_k \rightarrow w_k y_k$

$$\beta = (WX^T WX)^{-1} (WX^T WY)$$



Locally weighted regression is very flexible and fast to train.

Input X matrix of inputs: $X[k][i]$ = i'th component of k'th input point.

Input Y matrix of outputs: $Y[k]$ = k'th output value.

Input xq = query input. Input kwidth.

WXTWX = empty (D+1) x (D+1) matrix

WXTWY = empty (D+1) x 1 matrix

for (k = 1 ; k <= N ; k = k + 1)

/* Compute weight of kth point */

wk = weight_function(distance xq , X[k]) / kwidth)

/* Add to (WX) ^T (WX) matrix */

for (i = 0 ; i <= D ; i = i + 1)

for (j = 0 ; j <= D ; j = j + 1)

if (i == 0) xki = 1 else xki = X[k][i]

if (j == 0) xkj = 1 else xkj = X[k][j]

WXTWX [i][j] = WXTWX [i][j] + wk * wk * xki * xkj

/* Add to (WX) ^T (WY) vector */

for (i = 0 ; i <= D ; i = i + 1)

if (i == 0) xki = 1 else xki = X[k][i]

WXTWY [i] = WXTWY [i] + wk * wk * xki * Y[k]

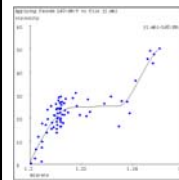
/* Compute the local beta. Call your favorite linear equation solver. Recommend Cholesky

Decomposition for speed. Recommend Singular Val Decomposition for Robustness. */

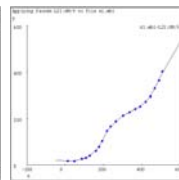
beta = (WXTWX)^-1 (WXTWY)

ypredict = beta[0] + beta[1]*xq[1] + beta[2]*xq[2] + ... beta[D]*xq[D]

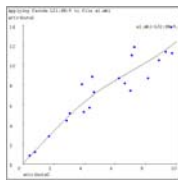
LWR on our test cases



KW = 1/16 of x-axis width.



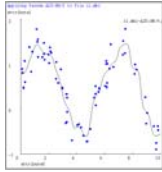
KW = 1/32 of x-axis width.



KW = 1/8 of x-axis width.

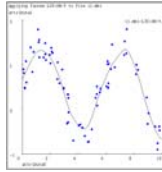
Nicer and smoother, but even now, are the bumps justified, or is this overfitting?

Locally weighted Polynomial regression



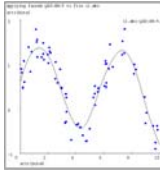
Kernel Regression
Kernel width K_W at optimal level.

$$KW = 1/100 \text{ x-axis}$$



LW Linear Regression
Kernel width K_W at optimal level.

$$KW = 1/40 \text{ x-axis}$$



LW Quadratic Regression
Kernel width K_W at optimal level.

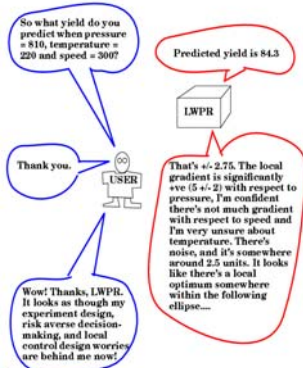
$$KW = 1/15 \text{ x-axis}$$

Local quadratic regression is easy: just add quadratic terms to the $WXTWX$ matrix. As the regression degree increases, the kernel width can increase without introducing bias.

Nice things about LWPR

1. It naturally extends to multiple inputs / multiple outputs.
2. There's no interference: new data doesn't cause unlearning of old parts of the model.
3. Incorporating new data is very cheap!
4. Handles functions where the output is linear in some variables and non-linear in others very naturally.
5. Leave-One-Out cross-validation is cheap
6. **You don't just get a prediction. You can get a full regression analysis!**

LWPR Regression Analysis



LWPR: The Dark Side

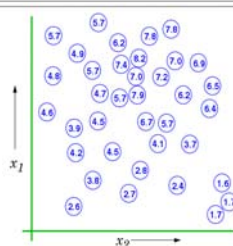
Cost of each prediction grows linearly with the number of datapoints.

With N datapoints, and with M polynomial terms...

- training is FREE
- each individual prediction costs $O(NM^2 + M^3)$

What can we do????

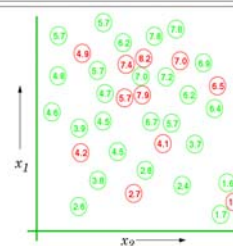
Idea One: Editing



Here we have $N=36$ examples, each with two inputs (x_1, x_2) and one output.

Instead of retaining and searching over all N , we could choose to keep only a few "representative" points called **prototypes**.

Idea One: Editing



- Examples: [Kibler and Aha, 1988; Skalak, 1994].
- Other editing tricks: maintain averages in the prototypes [Aha, 1992].

BUT: Lose accuracy. Lose ability to give different queries different kernel widths. Lose cheap L-O-O. Hard to estimate noise.

Idea 2: cache output surface

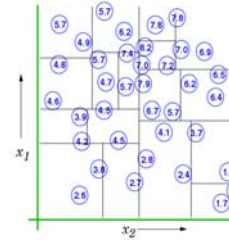
- First, decide which instance-based algorithm you'll use, fix the kernel width etc.
- Caching step:** perform lots of (slow) predictions but save the results in some interpolative table.
 - Simple (but expensive):** cache in a multidimensional array.
 - Fancy (still expensive):** cache in a kd -tree and predict with multilinear interpolation [Grosse, 1992].
 - Simple, cheaper, but discontinuous:** cache linear models inside leaves of a tree. e.g. [Quinlan, 1993].

BUT: Must commit to a kernel width during caching. Lose regression analysis. Caching step can be very very expensive. Incrementally adding points can be hard.

25

Idea 3: Range Searching

Embed your data in the leaves of a kd -tree [Friedman et al, 1977]

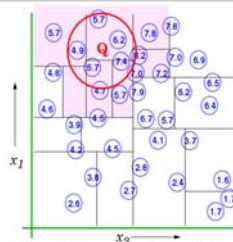


If you have a small kernel width, then beyond a certain distance from your query, all weights will be very close to zero.

You can quickly search the tree to get the set of points with non-zero-weight.

26

Idea 3: Range Searching



PRO: For small kernel widths much cheaper than linear search. Different queries can have different kernels. And you can do full regression analysis.

CON: For large kernel widths or very dense datasets, you're still in trouble.

27

This talk: Multiresolution Approach

What do we really need in order to do a locally weighted regression query?

$$\beta = (\sum_h w_h^2 \mathbf{x}_h^T \mathbf{x}_h)^{-1} (\sum_h w_h^2 \mathbf{x}_h^T \mathbf{y}_h)$$

Let's look at how we can build the

$(\sum_h w_h^2 \mathbf{x}_h^T \mathbf{x}_h)$ matrix...

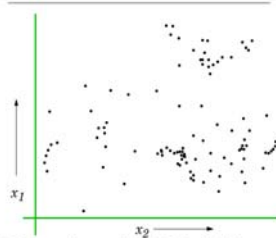
$$\begin{bmatrix} \sum_h w_h^2 x_{h1}^2 & \sum_h w_h^2 x_{h1} x_{h2} & \dots & \sum_h w_h^2 x_{h1} x_{hN} \\ \sum_h w_h^2 x_{h1} x_{h2} & \sum_h w_h^2 x_{h2}^2 & \dots & \sum_h w_h^2 x_{h2} x_{hN} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_h w_h^2 x_{h1} x_{hN} & \sum_h w_h^2 x_{h2} x_{hN} & \dots & \sum_h w_h^2 x_{hN}^2 \end{bmatrix}$$

The slow way:

```
wwxtx = MxM matrix of zeros.
for ( k = 1 ; k <= N ; k := k + 1 )
    wwxtx := wwxtx + w_h^2 x_h^T x_h
```

28

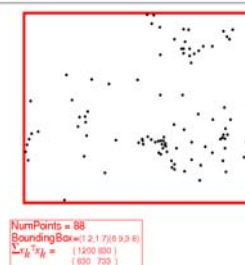
A multiresolution structure



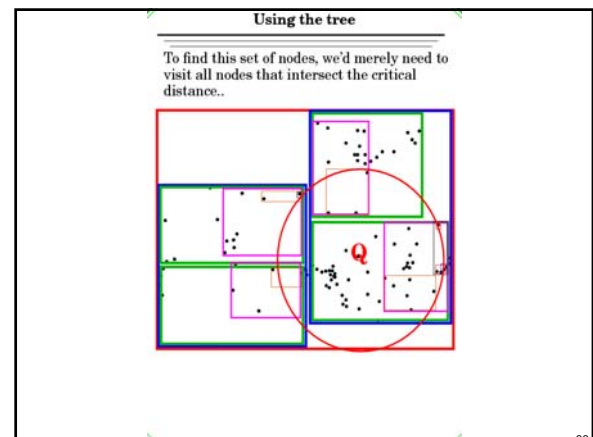
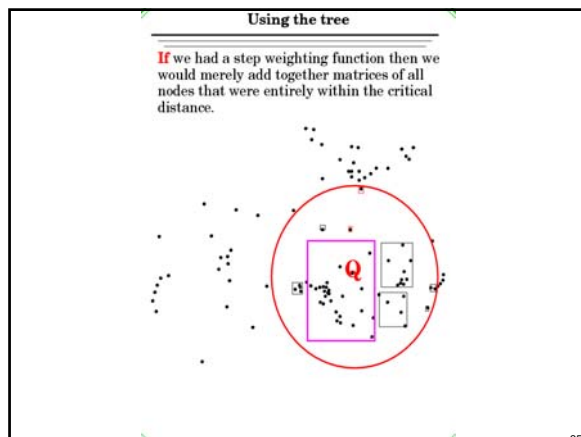
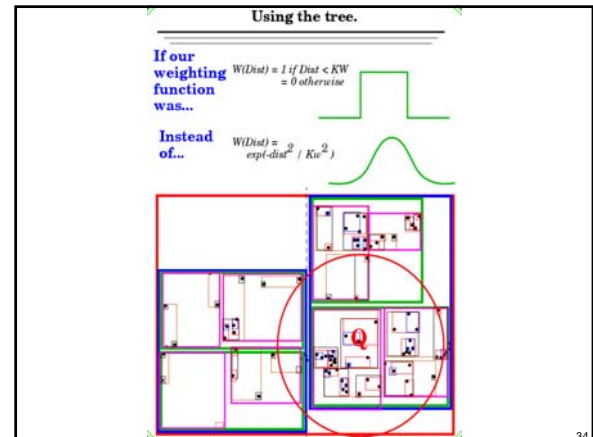
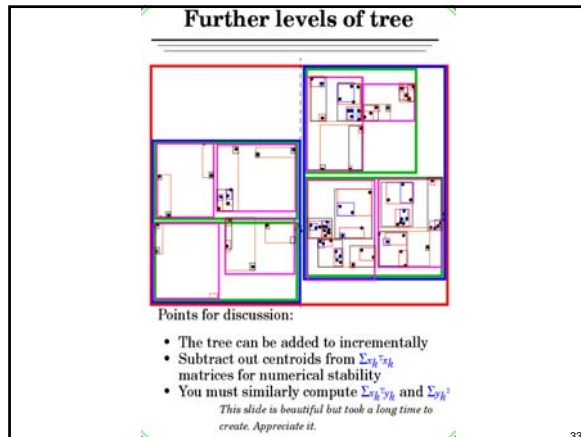
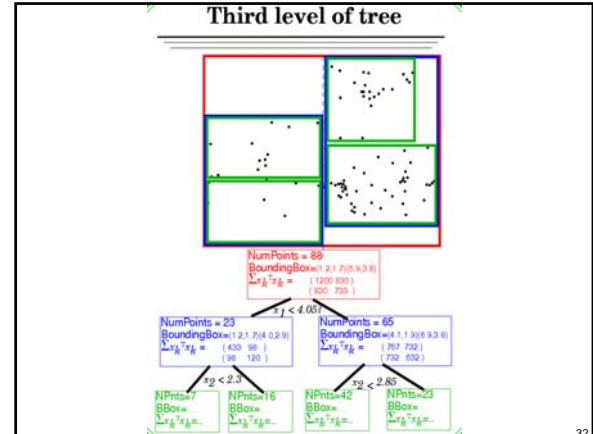
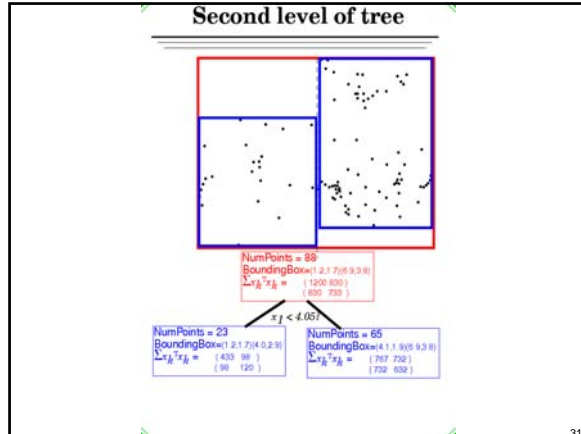
Take your data and build a kind of kd -tree structure with extra cached information...

29

Top level of tree



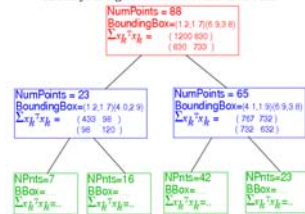
30



With a Gaussian weight fn...

Recurse down the tree.

- At each node you must compute the weighted sum $\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k$ of all points below in the tree.
- If node has only one point \mathbf{x}_j , return $w_j^2 \mathbf{x}_j^T \mathbf{x}_j$
- Else compute $\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k$ of your left subtree and in your right subtree. Return their sum.

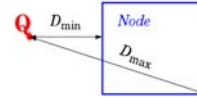


- But **SOMETIMES** you may discover that you can avoid searching a node...

37

Cutoff Criterion

You're about to compute $\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k$ for a node.



You know that all weights inside the node must lie in the range

$$w_{\min} = \text{weight}(D_{\max}) \leq w_k \leq \text{weight}(D_{\min}) = w_{\max}$$

So if w_{\max} and w_{\min} are very similar, then we know that

$$\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k = (w_{\max} + w_{\min}) / 2 * \sum_k \mathbf{x}_k^T \mathbf{x}_k$$

That information is inside the current node.

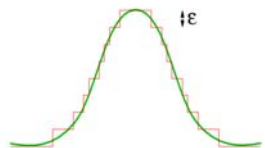
w_{\max} and w_{\min} are very similar when (1) the weight function is flat or (2) the node is small.

38

Cutoff Rule

To compute $\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k$ for a non-leaf node:

- If $w_{\max} - w_{\min} < \epsilon$ then return $(w_{\max} + w_{\min}) / 2 * \sum_k \mathbf{x}_k^T \mathbf{x}_k$
- Else recursively compute $\sum_k w_k^2 \mathbf{x}_k^T \mathbf{x}_k$ for left and right children and add them.



This is equivalent to using some approximate weighting function that never deviates from desired weighting function by more than ϵ

39

One final little detail...

- $\sum_k w_k^2$ can get as small as $O(\epsilon)$ if we are far from data.
- So approximations of $O(\epsilon)$ can become dangerously significant.

Thus we alter our criterion to:

$$\text{Cutoff if } w_{\max} - w_{\min} < \tau W_{\text{TOTAL}}$$

where τ is a small constant and where W_{TOTAL} is the total sum of squared weights of all datapoints in this query.

- W_{TOTAL} is not known in advance, but has a lower bound W_{SOFAR} . So in actuality...

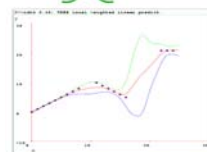
$$\text{Cutoff if } w_{\max} - w_{\min} < \tau W_{\text{SOFAR}}$$

40

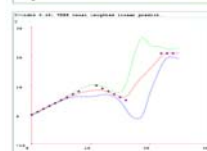
Slight Approximation

Examine the difference for locally weighted linear regression.

$KW = 4$



The linear, exact, method.



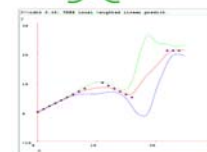
Multires, using $\tau = 10^{-7}$
..it's identical

41

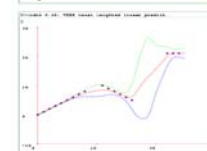
Medium Approximation

Examine the difference for locally weighted linear regression.

$KW = 4$



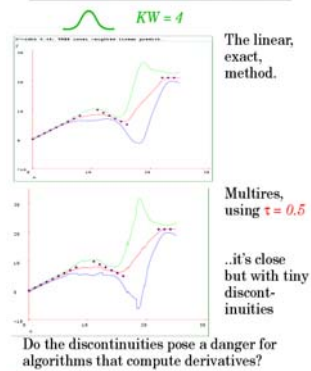
The linear, exact, method.



Multires, using $\tau = 0.05$
..it's indistinguishable

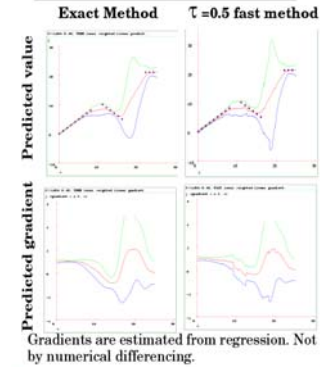
42

Gross Approximation



43

Gradient Estimation



44

Speed and Accuracy

Using locally weighted regression with kernel width 1/32 in unit-scaled coordinates.

	Expense	Error
30x31x301	Regular 12.25	0.22
1401x91	Regular 32.95	0.28
13 inputs	Time 22.23	0.28
170 datapoints	Approx 18.93	0.28
StdError 0.43	Fast 19.37	0.28
30x31x301	Regular 34.85	0.63
1401x91	Regular 33.43	0.63
3 inputs	Time 22.33	0.63
170 datapoints	Approx 4.43	0.63
StdError 2.2140	Fast 0.80	0.62
30x31x301	Regular 115.37	11.93
1401x91	Regular 151.30	11.93
5 inputs	Time 32.37	11.93
2144 datapoints	Approx 3.44	15.15
StdError 224.07	Fast 1.10	25.40
30x31x301	Regular 364.08	1.93
1401x91	Regular 469.00	1.93
10 inputs	Time 469.00	1.93
4077 datapoints	Approx 3.50	1.67
StdError 2.46	Fast 1.70	1.71
30x31x301	Regular 70.10	1.93
1401x91	Regular 10.10	1.93
3 inputs	Time 14.10	1.93
292 datapoints	Approx 11.45	1.93
StdError 0.9153	Fast 0.20	1.93
30x31x301	Regular 172.00	0.03
1401x91	Regular 176.10	0.03
3 inputs	Time 18.40	0.03
599 datapoints	Approx 11.82	0.03
StdError 0.3	Fast 0.21	0.02

Regular: standard linear method. Regular: Linear method, but ignores very weights. Time: Multires with $\tau = 1e-7$. Approx: $\tau = 0.05$. Fast: $\tau = 0.5$. Expense is seconds for 100 predictions.

45

Speed and Accuracy

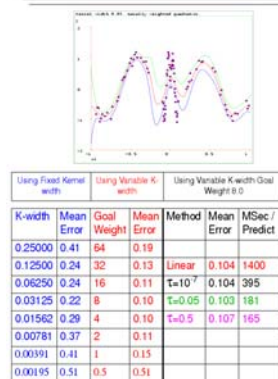
Using dataset-specific LWPR model selected by cross-validation.

	Expense	Error
30x31x301	Regular 37.86	0.22
1401x91	Regular 25.84	0.22
13 inputs	Time 18.32	0.22
170 datapoints	Approx 13.82	0.22
StdError 0.43	Fast 0.50	0.24
30x31x301	Regular 34.85	0.63
1401x91	Regular 33.43	0.63
3 inputs	Time 25.43	0.63
170 datapoints	Approx 8.17	0.63
StdError 2.2140	Fast 1.20	0.62
30x31x301	Regular 116.18	6.12
1401x91	Regular 151.30	6.12
5 inputs	Time 107.29	6.12
2144 datapoints	Approx 25.75	6.03
StdError 224.07	Fast 1.50	7.50
30x31x301	Regular 358.90	1.93
1401x91	Regular 469.00	1.93
10 inputs	Time 469.00	1.93
4077 datapoints	Approx 2.35	1.53
StdError 2.46	Fast 1.40	1.54
30x31x301	Regular 10.79	1.93
1401x91	Regular 10.10	1.93
3 inputs	Time 8.41	1.93
292 datapoints	Approx 1.70	1.94
StdError 0.9153	Fast 1.20	1.93
30x31x301	Regular 14.06	0.03
1401x91	Regular 14.96	0.03
3 inputs	Time 2.20	0.03
599 datapoints	Approx 2.20	0.03
StdError 0.3	Fast 0.50	0.02

Regular: standard linear method. Regular: Linear method, but ignores very weights. Time: Multires with $\tau = 1e-7$. Approx: $\tau = 0.05$. Fast: $\tau = 0.5$. Expense is seconds for 100 predictions.

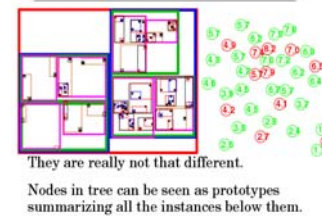
46

Variable kernel-widths



47

Multiresolution versus Prototypes



48

Comments

- Multiresolution trees make LWPR much faster without significant sacrifices.
- But predictions are still not free!
- This work extends [Deng+Moore 1995]'s kernel regression work.
- Similar multiresolution tricks are used in computer graphics and galaxy simulations
- Ongoing work plays similar tricks for local classification algorithms:
e.g. Locally weighted logistic regression [Deng+Moore, 1997]
- Note that with large datasets, the leaves of the tree are rarely visited...

Idea: Place the bottom levels of the tree in secondary memory (disk array) and maybe cope with billions of examples?