

Regression

3. Locally Weighted Regression Methods

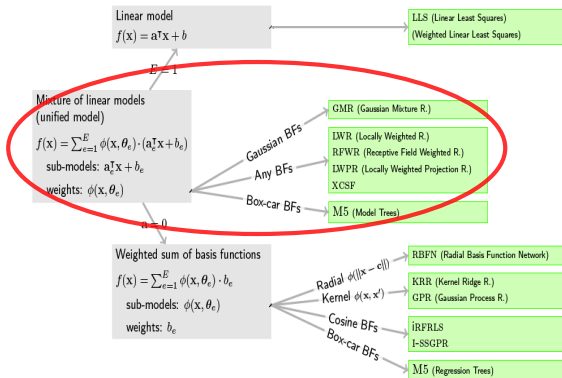
Olivier Sigaud

Sorbonne Université
<http://people.isir.upmc.fr/sigaud>



Locally Weighted Regression

Reminder: Outline of methods



- Multiple local and weighted least square regressions (shown with LWR)



Stulp, F. and Sigaud, O. (2015) Many regression algorithms, one unified model: A review. *Neural Networks*, 69:60–79

Locally Weighted Regression

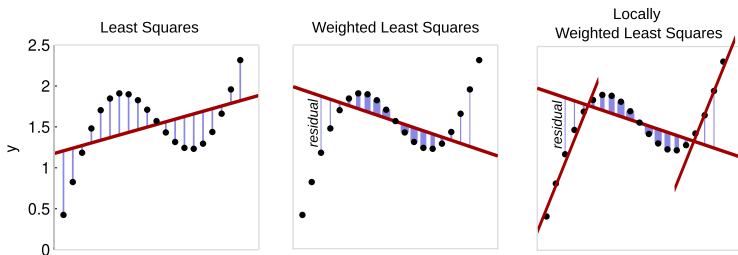


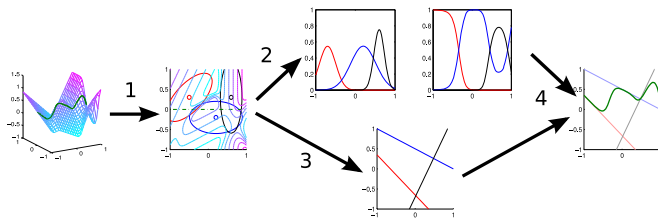
Figure: The thickness of the lines indicates the weights.

- ▶ General idea:
 - ▶ Split the function domain into linear parts
 - ▶ Give more weight to the centers
- ▶ Local linear models are tuned with Least Squares
- ▶ The importance of datapoints is represented by a Gaussian function



William S Cleveland and Susan J Devlin (1988) Locally weighted regression: an approach to regression analysis by local fitting.
Journal of the American statistical association, 83(403):596–610.

Locally Weighted Regression: Processes



1. Define the split into regions (**receptive fields**): by hand, data-driven, evolutionary
2. Determine relative importance of domains
3. Find linear models in the regions
4. Combine all linear models



Atkeson, C. (1991) Using locally weighted regression for robot learning. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, vol. 2, pp. 958–963

Combining linear models

- ▶ There are E features, or receptive fields (RF)
- ▶ Each RF is defined as a Gaussian $\phi(\mathbf{x}, \boldsymbol{\theta}_i) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x}-\boldsymbol{\mu}_i)}$ with $\boldsymbol{\theta}_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$
- ▶ To each RF is associated a local linear model

$$\Psi_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} + b_i$$

or $\Psi_i(\mathbf{x}) = w(\mathbf{x})^T \boldsymbol{\theta}_i$ and $w(\mathbf{x}) = (\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_d \ 1)^T$.

- ▶ Gaussians tell you how much each local linear model contributes to the output

$$y = \frac{\sum_{i=1}^E \phi(\mathbf{x}, \boldsymbol{\theta}_i) \Psi_i(\mathbf{x})}{\sum_{i=1}^E \phi(\mathbf{x}, \boldsymbol{\theta}_i)}$$

Batch learning (1)

- ▶ Consider a batch of N $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{1 \leq i \leq N}$ data elements.
- ▶ We want to find

$$y = f(\mathbf{x}) = \frac{\sum_{i=1}^E \phi(\mathbf{x}, \boldsymbol{\theta}_i) \Psi_i(\mathbf{x})}{\sum_{i=1}^E \phi(\mathbf{x}, \boldsymbol{\theta}_i)}$$

- ▶ Each local model should minimize the following locally weighted error:

$$\begin{aligned} \epsilon_i(\boldsymbol{\theta}_i) &= \frac{1}{2N} \sum_{j=1}^N \phi(\mathbf{x}^{(j)}, \boldsymbol{\theta}_i) \left(y^{(j)} - \Psi_i(\mathbf{x}^{(j)}) \right)^2 \\ &= \frac{1}{2N} \sum_{j=1}^N \phi(\mathbf{x}^{(j)}, \boldsymbol{\theta}_i) \left(y^{(j)} - w(\mathbf{x}^{(j)})^\top \boldsymbol{\theta}_i \right)^2. \end{aligned}$$

Batch learning (2)

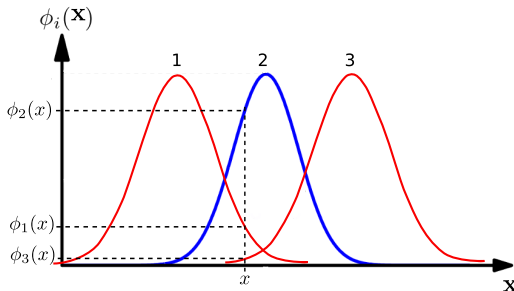
- ▶ As with the least squares method, the minimum is reached when the gradient is null:

$$-\frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}^{(j)}, \boldsymbol{\theta}_i) w(\mathbf{x}^{(j)}) \left(y^{(j)} - w(\mathbf{x}^{(j)})^\top \boldsymbol{\theta}_i \right) = 0.$$

- ▶ Therefore, we pose $\boldsymbol{\theta}_i = A_i^\# b_i$, with:

$$A_i = \sum_{j=1}^N \phi(\mathbf{x}^{(j)}, \boldsymbol{\theta}_i) w(\mathbf{x}^{(j)}) w(\mathbf{x}^{(j)})^\top$$
$$b_i = \sum_{j=1}^N \phi(\mathbf{x}^{(j)}, \boldsymbol{\theta}_i) w(\mathbf{x}^{(j)}) y^{(j)}.$$

LWR: algorithm

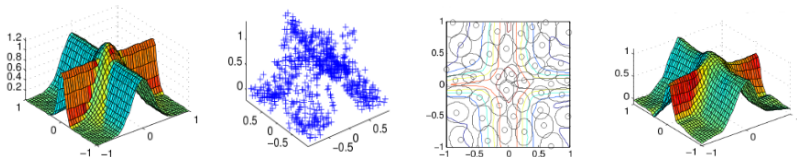


- ▶ LWR is a family of methods, but also the simplest algorithm of the family
- ▶ Gaussian features are evenly placed
- ▶ Adds receptive fields around non covered datapoints, moves the shape, cannot remove them
- ▶ Good along local trajectories



Atkeson, C. G., Moore, A. W., and Schaal, S. (1997) Locally weighted learning. *Lazy learning*, pages 11–73

LWPR: general goal

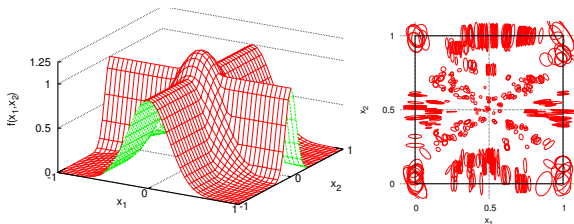


- ▶ Non-linear function approximation in very large spaces
- ▶ Using PLS to project linear models in a smaller space
- ▶ Adds receptive fields around non covered datapoints, moves the shape, cannot remove them
- ▶ Good along local trajectories



Schaal, S., Atkeson, C. G., and Vijayakumar, S. (2002). Scalable techniques from nonparametric statistics for real time robot learning. *Applied Intelligence*, 17(1):49–60.

XCSF: overview



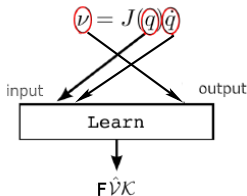
- ▶ XCSF is a Learning Classifier System [Holland, 1975]
- ▶ Linear models weighted by Gaussian functions (similar to LWPR)
- ▶ Linear models are updated using RLS
- ▶ Gaussian functions adaptation: Σ_i^{-1} and c_i are updated using a GA
- ▶ Condensation: reduce population to generalize better



Wilson, S. W. (2001) Function approximation with a classifier system. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2001)*, pages 974–981, San Francisco, California, USA. Morgan Kaufmann

XCSF: key feature

- ▶ Distinguish the space of linear models (prediction space) and the space of weights (condition space)
- ▶ Forward kinematics: $\dot{\xi} = F_{\theta}(\mathbf{q}, \dot{\mathbf{q}})$ $\dot{\xi} = J(\mathbf{q}) \dot{\mathbf{q}}$
- ▶ Forward dynamics: $\ddot{\mathbf{q}} = G_{\theta}(\mathbf{q}, \dot{\mathbf{q}}, \Gamma)$ $\ddot{\mathbf{q}} = A(\mathbf{q})^{-1} (\Gamma - \mathbf{n}(\mathbf{q}, \dot{\mathbf{q}}))$

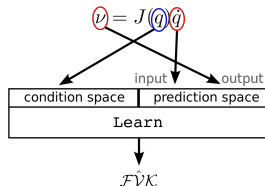


Forward kinematics with LWPR

- ▶ Example: learning forward kinematics ($\mathbf{x} = \langle \mathbf{q}, \dot{\mathbf{q}} \rangle$)
- ▶ LWPR: $\hat{f}(\mathbf{x}) = \sum_{i=1}^E \phi((\mathbf{q}, \dot{\mathbf{q}}), \theta_i) \cdot (b_i + \mathbf{a}_i^T(\mathbf{q}, \dot{\mathbf{q}}))$
- ▶ XCSF: $f(\mathbf{x}) = \sum_{i=1}^E \phi(\mathbf{q}, \theta_i) \cdot (b_i + \mathbf{a}_i^T \dot{\mathbf{q}})$

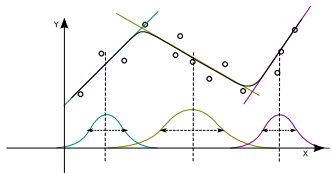


Butz, M., Pedersen, G., and Stalsh, P. (2009) Learning sensorimotor control structures with XCSF: redundancy exploitation and dynamic control. In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 1171–1178. ACM



Forward kinematics with XCSF

GMR



$$y = \sum_{k=1}^K h_k(\mathbf{x})(\mu_{k,Y} + \Sigma_{k,YX} \Sigma_{k,Y}^{-1}(\mathbf{x} - \mu_{k,X}))$$

With

$$\mu_k = [\mu_{k,X}^T, \mu_{k,Y}^T]^T \text{ and } \Sigma_k = \begin{pmatrix} \Sigma_{k,X} & \Sigma_{k,XY} \\ \Sigma_{k,YX} & \Sigma_{k,YY} \end{pmatrix}$$

- ▶ From input-output manifold to input-output function
- ▶ Same representation as the others, using $\theta^T = \Sigma_{i,YX} \Sigma_{i,Y}^{-1}$ and $b_i = \mu_{i,Y} - \Sigma_{i,YX} \Sigma_{i,X}^{-1} \mu_{i,X}$
- ▶ We get

$$\tilde{y} = \sum_{i=1}^E \frac{\pi_i \phi(\mathbf{x}, \theta_i)}{\sum_{l=1}^E \pi_l \phi(\mathbf{x}, \theta_l)} (\theta^T \mathbf{x} + b_i),$$

- ▶ Same as usual + scaling with the priors $\pi_i \rightarrow \pi_i = 1$ in standard model.
- ▶ Incorporates Bayesian variance estimation \rightarrow The richest representation



Hersch, M., Guenter, F., Calinon, S., & Billard, A. (2008) "Dynamical system modulation for robot learning via kinesthetic demonstrations." *IEEE Transactions on Robotics*, 24(6), 1463–1467

LWR methods: main features

Algo	LWR	LWPR	GMR	XCSF
Number of RFs	fixed	growing	fixed	adaptive
Position of RFs	fixed	fixed	adaptive	adaptive
Size of RFs	fixed	adaptive	adaptive	adaptive

- ▶ The main differences are in meta-parameter tuning
- ▶ Fewer hyperparameters is better, but less flexibility



Atkeson, C. (1991).

Using locally weighted regression for robot learning.

In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, volume 2, pages 958–963.



Atkeson, C. G., Moore, A. W., and Schaal, S. (1997).

Locally weighted learning.

Lazy learning, pages 11–73.



Butz, M., Pedersen, G., and Stalph, P. (2009).

Learning sensorimotor control structures with XCSF: redundancy exploitation and dynamic control.

In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 1171–1178. ACM.



Cleveland, W. S. and Devlin, S. J. (1988).

Locally weighted regression: an approach to regression analysis by local fitting.

Journal of the American statistical association, 83(403):596–610.



Hersch, M., Guenter, F., Calinon, S., and Billard, A. (2008).

Dynamical system modulation for robot learning via kinesthetic demonstrations.

IEEE Transactions on Robotics, 24(6):1463–1467.



Holland, J. H. (1975).

Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence.

University of Michigan Press, Ann Arbor, MI.



Schaal, S., Atkeson, C. G., and Vijayakumar, S. (2002).

Scalable techniques from nonparametric statistics for real time robot learning.

Applied Intelligence, 17(1):49–60.



Stulp, F. and Sigaud, O. (2015).

Many regression algorithms, one unified model: A review.

Neural Networks, 69:60–79.



Wilson, S. W. (2001).

Function approximation with a classifier system.

In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2001)*, pages 974–981, San Francisco, California, USA. Morgan Kaufmann.