Robust and Sparse

Combining Huber's loss with ℓ_1 constraint.

 $\hat{\boldsymbol{\alpha}}_{SparseHuber} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} J_{Huber}(\boldsymbol{\alpha})$ subject to $\|\boldsymbol{\alpha}\|_{1} \leq C$

- Solving quadratic programming problem is computationally rather demanding.
- Is it possible to make it faster?

I1 Loss

Quadratic term comes from Huber's loss. *ℓ*₁-loss is linear.





Linear Programming Learning ⁹¹

Combine ℓ_1 loss and ℓ_1 regularizer.



How to Obtain Solutions

Trick to avoid absolute value:

$$\|\boldsymbol{\alpha}\|_{1} = \min_{\boldsymbol{u} \in \mathbb{R}^{p}} \sum_{i=1}^{p} u_{i}$$

subject to $-\boldsymbol{u} \leq \boldsymbol{\alpha} \leq \boldsymbol{u}$,
Solution is given by
$$\min_{\boldsymbol{\alpha}, \boldsymbol{u} \in \mathbb{R}^{p}, \boldsymbol{h} \in \mathbb{R}^{n}} \left[\sum_{i=1}^{n} h_{i} + \lambda \sum_{i=1}^{p} u_{i} \right]$$

subject to $-\boldsymbol{h} \leq \boldsymbol{X} \boldsymbol{\alpha} - \boldsymbol{y} \leq \boldsymbol{h}$
 $-\boldsymbol{u} \leq \boldsymbol{\alpha} \leq \boldsymbol{u}$

Linearly Constrained Linear Programming Problem

93

Standard optimization softwares can solve the following form of linearly constrained linear programming problems.

Sparseness and Robustness ⁹⁴

	Sparse- ness	Robust- ness	Optimizat ion
ℓ_1 constrained LS	Yes	No	Quadratic
Huber's method	No	Yes	Quadratic
ℓ_1 constrained Huber	Yes	Yes	Quadratic
Linear programming	Yes	Yes	Linear

Non-Linear Models

A popular choice: Hierarchical models

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^p \alpha_i \varphi(\boldsymbol{x}; \boldsymbol{\beta}_i)$$
Basis function is parameterized by $\boldsymbol{\beta}_i$

Three-Layer Networks

Such a hierarchical model can be represented as a 3-layer network.



Sigmoidal Function

98

A typical basis function: Sigmoidal function

$$\varphi(\boldsymbol{x};\boldsymbol{\beta},b) = rac{1}{1 + \exp\left(-\langle \boldsymbol{x},\boldsymbol{\beta} \rangle - b\right)}$$



Perceptrons

- The behavior of the sigmoidal functions is similar to the neurons in the brain.
- For this reason, hierarchical models with sigmoidal functions are called artificial neural networks or perceptrons.
- Mathematically, 3-layer neural networks can approximate any continuous functions with arbitrary small error ("universal approximator").

Gaussian Radial Basis Function¹⁰⁰

Another popular basis function: Gaussian radial basis function



Non-Linear Learning for Non-Linear Models

102

Sigmoidal neural networksError back-propagation algorithm

Least-squares Learning $\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{p} \alpha_i \varphi(\boldsymbol{x}; \boldsymbol{\beta}_i)$ $\boldsymbol{w} = (\boldsymbol{\alpha}^{\top}, \boldsymbol{\beta}_1^{\top}, \boldsymbol{\beta}_2^{\top}, \dots, \boldsymbol{\beta}_p^{\top})^{\top}.$

103

Least-squares learning is often used for training hierarchical models.

$$\min_{\boldsymbol{w}} J_{LS}(\boldsymbol{w})$$
$$J_{LS}(\boldsymbol{w}) = \sum_{i=1}^{n} \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)^2$$

How to Obtain Solutions

104

No analytic solution is known.We usually use simple gradient search.



It converges to one of the local minima.

Error Back-Propagation

Efficient Calculation of Gradient for Sigmoidal Basis Functions

$$\frac{\partial J_{LS}}{\partial \alpha_j} = 2 \sum_{i=1}^n o_{i,j} \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)$$
$$\frac{\partial J_{LS}}{\partial \beta_j^{(k)}} = 2\alpha_j \sum_{i=1}^n o_{i,j} \left(1 - o_{i,j} \right) x_i^{(k)} \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)$$
$$\frac{\partial J_{LS}}{\partial b_j} = 2\alpha_j \sum_{i=1}^n o_{i,j} \left(1 - o_{i,j} \right) \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)$$

 $o_{i,j} = \varphi(\boldsymbol{x}_i; \boldsymbol{\beta}_j, b_j)$

105

Error Back-Propagation (cont.)¹⁰⁶

- When the output values of the network are calculated, the input points are propagated following the forward path.
- On the other hand, when the gradients are calculated, the output error $(\hat{f}(\boldsymbol{x}_i) y_i)$ is propagated backward.
- For this reason, this algorithm is called the error back-propagation.
- However, it is gradient search so global convergence is not guaranteed.

Stochastic Gradient Descent¹⁰⁷

- In the usual gradient method, all training examples are used at the same time.
- In practice, the following stochastic method would be computationally advantageous.
 - Randomly choose one of the training examples (say, $({m x}_i, y_i)$)
 - Update the parameter vector by

$$\widehat{\boldsymbol{w}}^{new} \longleftarrow \widehat{\boldsymbol{w}}^{old} - \varepsilon \nabla J_i(\widehat{\boldsymbol{w}}^{old})$$
$$J_i(\boldsymbol{w}) = (\widehat{f}(\boldsymbol{x}_i) - y_i)^2$$

• Repeat this procedure until convergence.