## Non-linear least Squares

unknown true model: $\quad y^* = f_{\beta^*}(x) + \varepsilon \qquad\qquad \varepsilon \sim N(0, \sigma^2)$

$f_{\beta^*}(x)$ : non-linear fct with true parameters $\beta^*$

[special case: $f_{\beta^*}(x) = x \cdot \beta^*$ (linear fct.) $\Rightarrow$ ordinary least squares]

goal: $\quad \hat{\beta} = \underset{\beta}{\arg\min} \; \|Y - f_\beta(x)\|^2$ , in general no analytic solution

many approximations:
- regression trees, regression forests
- Levenberg - Marquardt algorithm
- change of variables to transform non-linear probl. into a linear one (approximately)
- kernel methods
- Gaussian processes
- neural networks
- ...

FML { ... }

Advanced Machine Learning { ... }

## Regression Trees

work like density and decision trees: partition the feature space into bins using a tree, have a constant response for every bin :

$$\hat{Y_i} = \hat{f}(x_i) = \sum_{\ell=1}^{L} f_\ell \; \mathbb{1}[x_i \in bin_\ell]$$

sum over all bins ↑ constant bin response ↑ fct to find the bin for $x_i$

find optimal response in bin$_\ell$ via least squares loss (assume bins are already fixed):

$$\hat{f_\ell} = \int_{x \in bin_\ell} (f_\ell - f(x))^2 \, dx \qquad \hat{f}(x) = \sum_{\ell=1}^{L} f_\ell \, \mathbb{1}[x_i \in bin_\ell]$$

$$\{\hat{f_\ell}\} = \arg\min_{\{f_\ell\}} \int (\hat{f}(x) - f(x))^2 \, p(x) \, dx$$

$$= \arg\min_{\{f_\ell\}} \underbrace{\int f(x)^2 p(x) \, dx}_{\substack{\text{independent of model} \\ \rightsquigarrow \text{drop}}} - 2\int \hat{f}(x) f(x) \, p(x) \, dx + \int \hat{f}(x)^2 p(x) \, dx$$

$$= \arg\min_{\{f_\ell\}} -2 \sum_{\ell=1}^{L} f_\ell \int_{x \in bin_\ell} f(x) \, p(x) \, dx + \sum_{\ell} f_\ell^2 \int_{x \in bin_\ell} p(x) \, dx$$

$$\frac{\partial Loss}{\partial f_\ell} = -2 \underbrace{\int_{x \in bin_\ell} f(x) \, p(x) \, dx}_{\substack{\mathbb{E}[f(x)] \\ bin_\ell}} + 2 f_\ell \underbrace{\int_{x \in bin_\ell} p(x) \, dx}_{\text{total density in } bin_\ell} \overset{!}{=} 0$$

$$\approx \left(\frac{1}{N_\ell} \sum_{i: x_i \in bin_\ell} f(x_i)\right) \cdot \frac{N_\ell}{N} + f_\ell \frac{N_\ell}{N} \overset{!}{=} 0 \qquad \overbrace{\hat{f_\ell} = \frac{1}{N_\ell} \sum_{i: x_i \in bin_\ell} f(x_i)}^{=\, y_i}$$

$$\boxed{\hat{f_\ell} = \frac{1}{N_\ell} \sum_{i \in bin_\ell} y_i}$$

regression trees: piece-wise constant approximation of a nonlinear function

truth: $Y = f(X) + \varepsilon$    approx: $y = \hat{f}(x) = \sum_{\ell=1}^{L} \hat{f}_\ell \, \mathbb{1}[X \in bin_\ell]$

minimize the squared loss

$$\{\hat{f}_\ell\}_{\ell=1}^{L} = \underset{\{f_\ell\}}{\arg\min} \; \mathbb{E}_{p(x)}\left[ (f(x) - \hat{f}(x))^2 \right] = \underset{\{f_\ell\}}{\arg\min} \int (f(x) - \hat{f}(x))^2 \, p(x) \, dx$$

$$= \underset{\{f_\ell\}}{\arg\min} \; -2 \sum_{\ell=1}^{L} f_\ell \underbrace{\int_{x \in bin_\ell} f(x) \, p(x) \, dx}_{\displaystyle \sum_{i: x_i \in bin_\ell} Y_i (= f(x_i)) \frac{1}{N}} + \sum_{\ell=1}^{L} f_\ell^2 \underbrace{\int_{x \in bin_\ell} 1 \cdot p(x) \, dx}_{\displaystyle \sum_{i: x_i \in bin_\ell} 1 \cdot \frac{1}{N}}$$

$$= \underset{\{f_\ell\}}{\arg\min} \; -2 \sum_{\ell=1}^{L} f_\ell \cdot \frac{1}{N} \sum_{i: x_i \in bin_\ell} Y_i + \sum_{\ell=1}^{L} f_\ell^2 \frac{N_\ell}{N}$$

$$\frac{\partial \, loss}{\partial f_\ell} = -2 \frac{1}{N} \sum_{i: x_i \in bin_\ell} Y_i + 2 f_\ell \frac{N_\ell}{N} \overset{!}{=} 0$$

$$\boxed{\hat{f}_\ell = \frac{1}{N_\ell} \sum_{i: x_i \in bin_\ell} Y_i}$$

**training alg.:** exactly like decision & density tree, just with different criterion for selecting the optimal split! in each node to be split, iterate over all candidate splits and compute the squared error: $f_{left} = \frac{1}{N_{left}} \sum_{i: x_i \in left} Y_i$  $f_{right} = \frac{1}{N_{right}} \sum_{i: x_i \in right} Y_i$

execute split with minimal R: $R = \sum_{i: x_i \in left} (Y_i - f_{left})^2 + \sum_{i: x_i \in right} (Y_i - f_{right})^2$

regression forest: to avoid overfitting, train an ensemble of regression trees with
the usual tricks ( each tree is trained with a bootstrap sample of the TS,
in each node a random subset of the features is considered for splitting)
and return the average response of all trees

## Levenberg ~ Marquardt - Algorithm

setting here: $\qquad y = f_\beta (x) + \epsilon \qquad$ with $f_\beta (x)$ a non-linear function with params
$\qquad$ parameters $\beta$, functional form (formula) of $f_\beta$ is
$\qquad$ known, only the true parameters $\beta^*$ are unknown

[notice: linear models $\cancel{x_\beta}\; f_\beta (x) = X \cdot \beta$ are a special case ]

difficulty: - for $f_\beta (x)$ non-linear, there is usually no analytic solution for $\beta$

$\Rightarrow$ reduce to a sequence of linear problems via Taylor series,
$\qquad$ iteratively solve non-linear optimization by solving a least-squares prob. in every iteration

least-squares objective: $\qquad \hat\beta = \arg\min_\beta \sum_{i=1}^{N} (y_i - f_\beta (x_i))^2$

- iterative Gauss - Newton alg. want to minimize $\mathcal{L}(\beta)$
  - given current guess $\beta^{(t-1)}$, expand $\mathcal{L}(\beta)$ into Taylor series at $\beta^{(t-1)}$

$$\underbrace{\mathcal{L}(\beta^{(t-1)} + \Delta)}_{= \beta^{(t)}} \approx \mathcal{L}(\beta^{(t-1)}) + \underbrace{\frac{\partial \mathcal{L}}{\partial \beta}\bigg|_{\beta^{(t-1)}}}_{\text{"Jacobian of the loss"}} \cdot \Delta \qquad + \text{higher orders (ignored)}$$

- least-squares minimization of w.r.t. improvement $\Delta$

$$\hat{\Delta} = \underset{\Delta}{\arg\min} \left( \underbrace{L(\beta^{(t-1)})}_{\tilde{Y}} + \underbrace{\frac{dL}{d\beta}\Big|_{\beta^{(t-1)}}}_{-\tilde{X}} \cdot \Delta \right)^2$$

auxilliary variables:    $\tilde{Y}$    (vector of length $N$)    $-\tilde{X}$    (matrix of derivatives $N \times D$)

$$\hat{\Delta} = \underset{\Delta}{\arg\min} \left( \tilde{Y} - \tilde{X} \cdot \Delta \right)^2 \qquad \hat{\Delta} = (\tilde{X}^T \tilde{X})^{-1} \cdot \tilde{X}^T \cdot \tilde{Y}$$

solution via pseudo-inverse

not yet optimal, because it tends to overfit $\Rightarrow$ $L_2$ regularization

$$\hat{\Delta} = \underset{\Delta}{\arg\min} \left( \tilde{Y} - \tilde{X} \cdot \Delta \right)^2 + \tau \|\Delta\|^2$$

- in case of non-linear least squares :    ~~$L(x)$~~

$$L(\beta) = Y - f_\beta(x)$$

$$L(\beta^{(t-1)} + \Delta) = \underbrace{L(\beta^{(t-1)})}_{= Y - f_{\beta^{(t-1)}}(x)} - \underbrace{\frac{df_\beta(x)}{d\tau}\Big|_{\beta^{(t-1)}}}_{\tilde{X}} \cdot \Delta$$

$$= \tilde{Y} \qquad\qquad - \tilde{X} \cdot \Delta$$

optimize via ridge-regression :    $\hat{\Delta} = (\tilde{X}^T \tilde{X} + \tau \mathbb{I})^{-1} \tilde{X}^T \tilde{Y}$

$$\beta^{(t)} = \beta^{(t-1)} + \Delta$$

- why is it good to use ridge regression ?

◊ case 1: $\tilde{X}^T\tilde{X}$ has good condition $\Rightarrow$ $\tau \approx 0$ (no regularization needed)

$\Rightarrow$ $\hat{\Delta}$ is unbiased least-squares solution $\Rightarrow$ very fast convergence (few iterations)

case 2: $\hat{X}^T\tilde{X}$ has bad condition $\Rightarrow$ $\tau \gg 0$,

in the extreme, $\tilde{X}^T\tilde{X}$ is dominated by $\tau \mathbb{I}$

$(\tilde{X}^T\tilde{X} + \tau\mathbb{I})^{-1}\tilde{X}^T\tilde{y} \Rightarrow \frac{1}{\tau}\tilde{X}^T\tilde{y}$

$\Rightarrow$ we get a gradient descent alg. with learning rate $\frac{1}{2\tau}$

$$\frac{\partial(\tilde{y} - \tilde{X}\cdot\Delta)^2}{\partial\Delta}\bigg|_{\Delta=0} = -2\tilde{X}^T\tilde{y}$$

$\Rightarrow$ $\tau$ allows us to interpolate between fast Newton iterations and slower, but numerically stable, gradient descent

. ✦ Questions: – how to choose good $\tau$? $\Rightarrow$ self-regulating algorithm

– how to ensure that $\tau$ has the same effect on all features

$\left[\text{recall: in ridge regression, we standardized the features } \tilde{X}_j = \dfrac{X_j - \bar{X}_j}{\text{std}(X_j)}\right]$

standardization is now impossible, because $\tilde{X}$ is the Jacobian matrix $\left(\frac{\partial f}{\partial\beta}\right)$, which would become incorrect by scaling

$\Rightarrow$ Marquardt's trick: weight $\tau$ with $\|\tilde{X}_j\|^2 = \text{diag}(\tilde{X}^T\tilde{X})$

$\hat{\Delta} = (\tilde{X}^T\tilde{X} + \tau \cdot \text{diag}(\tilde{X}^T\tilde{X}))^{-1}\cdot\tilde{X}^T\tilde{y}$

(in practice: compute $S = \tilde{X}^T\tilde{X}$, then scale diagonal elements $S_{jj} \Leftarrow S_{jj}(1+\tau)$

## full Levenberg - Marquardt algorithm

(0) define initial guess $\beta^{(0)}$ (this is usually critical, bad guess $\Rightarrow$ bad optimum)

initial regularization strength $\tau^{(0)}$

(1) for $t = 1, \dots, T$ (or until convergence):

(a) compute $\hat{y} = Y - f_{\beta^{(t-1)}}(X)$, $\hat{X} = \dfrac{\partial f_\beta(x)}{\partial \beta}\Big|_{\beta^{(t-1)}}$

(b) solve $\left(\hat{X}^T\hat{X} + \tau^{(t-1)} \text{diag}(\hat{X}^T\hat{X})\right) \cdot \Delta = \hat{X}^T \cdot \hat{y}$ by some

linear solver $\Rightarrow \hat{\Delta}$

(c) compute $\hat{\beta} = \beta^{(t-1)} + \hat{\Delta}$ and residuals $(Y - f_{\hat{\beta}}(X))^2$

(d) if residual (squared loss) of $\hat{\beta}$ is smaller than of $\beta^{(t-1)}$

- accept $\hat{\beta}$ as our new guess $\beta^{(t)} = \hat{\beta}$
- reduce regularization strength $\tau^{(t)} = \tau^{(t-1)}/\sqrt{2}$

[hyper param]

otherwise

- reject the guess $\hat{\beta}$: $\beta^{(t)} = \beta^{(t-1)}$
- increase the regularization $\tau^{(t)} = \sqrt{2}\,\tau^{(t-1)}$

$$\text{diag}(A) = \begin{bmatrix} A_{11} & & & 0 \\ & A_{22} & & \\ & & \ddots & \\ 0 & & & A_{ND} \end{bmatrix}$$