

8.2.2 Maximum likelihood inference

Today

- Boosting (Adaboost : the first popular boosting algorithm)
- statistical interpretation of boosting (L2 Boosting)
- likelihood-based boosting (together with model-based boosting (gradient boosting), statistical boosting)

Y_i iid with density $p(y_i; \theta)$

e.g. $Y_i \sim \mathcal{N}(\mu, \sigma^2) \rightarrow p(y_i; \underbrace{\mu, \sigma^2}_{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\right\}$

$$L(\theta; y) = \prod_{i=1}^N p(y_i; \theta) \quad \text{likelihood}$$

$$\ell(\theta; y) = \sum_{i=1}^N \log p(y_i; \theta) \quad \text{log-likelihood}$$

$$\hat{\theta} = \arg\max_{\theta} \ell(\theta; y) = \arg\max_{\theta} L(\theta; y)$$

$$\ell_{\theta}(\theta; y) = \frac{\partial \ell(\theta; y)}{\partial \theta} \quad \text{score function}$$

$$\ell_{\theta\theta}(\theta; y) = \frac{\partial^2 \ell(\theta; y)}{\partial \theta \partial \theta^T} \quad \begin{cases} \rightarrow j(\theta) = -\ell_{\theta\theta}(\theta; y) & \text{observed information} \\ \leftarrow i(\theta) = E_{\theta}[j(\theta)] & \text{expected information} \end{cases}$$

If θ_0 is the true parameter

$$\hat{\theta} \xrightarrow{n \rightarrow \infty} \mathcal{N}(\theta_0; j(\theta_0)^{-1})$$

Estimation

$$\hat{\theta} \sim \mathcal{N}(\hat{\theta}; j(\hat{\theta})^{-1}) \quad \text{or} \quad \hat{\theta} \sim \mathcal{N}(\hat{\theta}; i(\hat{\theta})^{-1})$$

confidence intervals $\hat{\theta} \pm z_{1-\frac{\alpha}{2}} \sqrt{j^{-1}(\hat{\theta})}$

10 Boosting

Leo Breiman: "[boosting is] the best off-the-shelf classifier in the world"

- originally developed for classification;
- translated into the statistical world and use for all purposes (regression, ...)
- extended in its use;
- interpretable (GAM)

Starting challenge:

"Can a committee of blackheads somehow arrive at a highly reasoned decision, despite the weak judgement of the individual members?"

goal: obtain a good classifier

- apply "weak estimators", in the context of classification, classifiers which lead to a solution only slightly better than a random choice

idea: repeatedly apply a weak estimator to modifications of the data (iteratively)

at each iteration gives more weight to the misclassified observations

↓
AdaBoost

Consider a two class classification problem

$$Y_i \in \{-1, 1\}$$

X_i : the vector of inputs

AdaBoost algorithm

① initialize, weights are $(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}) = w^{[0]}$

② for m from 1 to m-stop (M)

a) fit the weak estimator $\hat{g}(x)$ to the weighted data;

b) compute the weighted in-sample misclassification rate

$$\text{err}^{[m]} = \frac{\sum_{i=1}^N w_i^{[m-1]} \mathbb{1}[y_i \neq \hat{g}^{[m]}(x_i)]}{\sum_{i=1}^N w_i^{[m-1]}}$$

c) compute $\alpha_m = \log\left(\frac{1 - \text{err}^{[m]}}{\text{err}^{[m]}}\right)$

α_m is used to weight the contribution of the $\hat{g}^{[m]}(x)$ to the final estimate (classification)

d) update the weights

$$\hat{w}_i = w^{[m-1]} \exp\left\{\alpha_m \underbrace{\mathbb{1}[y_i \neq \hat{g}^{[m]}(x_i)]}_{\text{reweight only misclassified observations}}\right\} \quad i=1, \dots, N$$

or

$$w_i^{[m]} = \frac{\hat{w}_i}{\sum_{i=1}^N \hat{w}_i}$$

③ compute final result

$$\hat{f}_{\text{AdaBoost}} = \text{sgn}\left(\sum_{m=1}^{m\text{-stop}} \alpha^{[m]} \hat{g}^{[m]}(x)\right)$$

Example

step [0] $w = (\frac{1}{10}, \frac{1}{10}, \dots, \frac{1}{10})$

step [1] $err = \frac{\sum_{i=1}^N \frac{1}{10} \mathbb{1}[y_i \neq \hat{g}_i]}{\sum_{i=1}^N \frac{1}{10}} = \frac{\frac{3}{10}}{\frac{10}{10}} = 0.3$

$$\alpha_1 = \log \frac{1-err}{err} = \log 0.7 - \log 0.3 \approx 0.84$$

$$\tilde{w} = \left(\underbrace{\exp(0.84) \frac{1}{10}}_{\text{miss classifier in the first iteration}}, 0.23, 0.23, 0.1, \dots, 0.1 \right)$$

$$w^{[1]} \approx (0.17, 0.17, 0.17, 0.07, \dots, 0.07)$$

Statistical view of Boosting

- functional gradient descent algorithm
- forward stagewise (additive) modelling

→ see AdaBoost as an iterative procedure to minimize a loss-function, in particular an exponential loss-function

$$L(y, f(x)) = \exp\{-y f(x)\}$$

Consider a generic step m

- the current classifier is $\hat{f}^{(m-1)} = \sum_{k=1}^{m-1} \alpha_k \hat{g}^{(k)}(x)$ 1...m-1 given in learning
(if the algorithm had stopped at the $m-1$ iteration)

- the goal is to find $(\alpha_m, g_m) = \arg \min_{\alpha, g} \sum_{i=1}^N \exp\{-y_i \sum_{k=1}^m \alpha_k g^{(k)}(x_i)\}$

$$\begin{aligned} \alpha_m, g_m &= \arg \min_{\alpha, g} \sum_{i=1}^N \exp\left\{-y_i \left(\underbrace{\sum_{k=1}^{m-1} \alpha_k \hat{g}^{(k)}(x_i)}_{\hat{f}^{(m-1)}} + \alpha g\right)\right\} \\ &= \arg \min_{\alpha, g} \sum_{i=1}^N w_i^{(m-1)} \exp\{-y_i \alpha g\} \quad (*) \end{aligned}$$

where $w_i = \exp\{-y_i \hat{f}^{(m-1)}(x_i)\}$, which do not depend neither on α nor on g (given from the previous iteration)

- two step procedure: first we minimize with respect to g

$$\begin{aligned} & \arg \min_g \sum_{i=1}^N w_i^{(m-1)} \exp(-y_i \alpha g) \quad \begin{array}{l} \text{if } y_i = g \rightarrow y_i g = 1 \\ \text{if } y_i \neq g \rightarrow y_i g = -1 \end{array} \\ &= \arg \min_g \left\{ \sum_{g=y} w_i^{(m-1)} e^{-\alpha} + \sum_{g \neq y} w_i^{(m-1)} e^{\alpha} \right\} \\ &= \arg \min_g \left\{ e^{-\alpha} \sum_{i=1}^N w_i^{(m-1)} + (e^{\alpha} - e^{-\alpha}) \sum_{g \neq y} w_i^{(m-1)} \right\} \\ &= \arg \min_g \left\{ \sum_{g=y} w_i^{(m-1)} e^{-\alpha} + \sum_{g \neq y} w_i^{(m-1)} e^{-\alpha} - \sum_{g \neq y} w_i^{(m-1)} e^{-\alpha} + \sum_{g \neq y} w_i^{(m-1)} e^{\alpha} \right\} \\ &= \arg \min_g \left\{ e^{-\alpha} \sum_{i=1}^N w_i^{(m-1)} + (e^{\alpha} - e^{-\alpha}) \sum_{i=1}^N w_i^{(m-1)} \mathbb{1}[g \neq y_i] \right\} \\ & \underline{g = \arg \min_g \left\{ \sum_{i=1}^N w_i^{(m-1)} \mathbb{1}[g \neq y_i] \right\}} \end{aligned}$$

$$= \arg \min_{\alpha, g} \sum_{i=1}^n w_i^{[n-1]} \exp\{-y_i \alpha g\}$$

$$\arg \min_{\alpha} \sum w_i^{[n-1]} \exp\{-y_i \alpha g\} \quad \frac{\partial \mathcal{L}}{\partial \alpha} = 0$$

$$y = g \rightarrow \sum_{g=y} w_i^{[n-1]} \exp\{-\alpha\}$$

$$y \neq g \rightarrow \sum_{g \neq y} w_i^{[n-1]} \exp\{\alpha\}$$

$$\frac{\partial \mathcal{L}}{\partial \alpha} = -\sum_{g=y} w_i^{[n-1]} \exp\{-\alpha\} + \sum_{g \neq y} w_i^{[n-1]} \exp\{\alpha\} = 0$$

multiply
both terms
for e^{α}

$$= -\sum_{g=y} w_i^{[n-1]} + \sum_{g \neq y} w_i^{[n-1]} \exp\{2\alpha\} = 0$$

$$e^{2\alpha} \sum_{g \neq y} w_i^{[n-1]} = \sum_{g=y} w_i^{[n-1]}$$

$$e^{2\alpha} = \frac{\sum_{i=1}^n w_i^{[n-1]} - \sum_{g=y} w_i^{[n-1]}}{\sum_{g \neq y} w_i^{[n-1]}}$$

$$\alpha = \frac{1}{2} \log \left(\frac{1 - \text{err}}{\text{err}} \right)$$

$$\text{where } \text{err} = \frac{\sum_{i=1}^n w_i^{[n-1]} \mathbb{1}(y_i \neq g)}{\sum_{i=1}^n w_i^{[n-1]}}$$

- $\hat{g}^{[n]}$ minimizer of the weighted misclassification
- $\alpha = \frac{1}{2} \log \left(\frac{1 - \text{err}}{\text{err}} \right)$

Our general classifier is updated as

$$\hat{f}^{[m]} = \hat{f}^{[m-1]} + \alpha_m \hat{g}^{[m]}$$

which causes the weights of the next iteration to be

$$w_i^{[m]} = w_i^{[m-1]} \cdot \exp\{-\alpha y_i \hat{g}_i^{[m]}\}$$

Since $-y_i \hat{g}_i^{[m]} = -\sum_{g_j=y_i} 1 + \sum_{g_j \neq y_i} (+1) = \sum_{g_j \neq y_i} 1 - \sum_{g_j=y_i} 1$

$$= 2 \sum_{g_j \neq y_i} 1 - 1$$

$$w_i^{[m]} = w_i^{[m-1]} \cdot \exp\{+\alpha (2 \sum_{g_j \neq y_i} 1 - 1)\}$$

$$= w_i^{[m-1]} \exp\{2\alpha \mathbb{1}(y_i \neq g_i) - \alpha\}$$

$$= \underbrace{w_i^{[m-1]}}_{\text{constant for each observation} \rightarrow \text{can be ignored}} e^{-\alpha} e^{2\alpha \mathbb{1}(y_i \neq g_i)}$$

2 β is α in the algorithm

→ AdaBoost minimizes the exponential loss criterion by a forward stagewise procedure

Note:

- this statistical view allow us to interpret the results of the procedure.

In particular, it can be showed that the minimizer of the exponential loss

$$f^*(x) = \arg \min_{f(x)} E_{Y|x} [e^{-y f(x)}] = \frac{1}{2} \log \frac{\Pr[Y=1|x]}{\Pr[Y=-1|x]}$$

alternatively: $\Pr[Y=1|x] = \frac{1}{1 + e^{-2f^*(x)}}$

$\frac{1}{2}$ the log-odds for $\Pr[Y=1|x]$ → Ex 10.2

• other loss-functions lead to the same minimizer, for example the negative log-likelihood

$$\hat{\pi} = \Pr[Y=1|x] = \frac{e^{f(x)}}{e^{f(x)} + e^{-f(x)}} = \frac{1}{1 + e^{-2f(x)}}$$

$$y' = \frac{y+1}{2} \in \{0,1\}$$

then

$$\ell(\pi) = y' \log \pi + (1-y') \log (1-\pi)$$

$$\Rightarrow -\ell(\pi) = \log(1 + e^{-2y' f(x)})$$

↙₂

