## Weighted Least Squares

- when the noise is different for each instance $i$
- recall: OLS assumes that noise variance $\sigma^2$ is equal for all $i$
- maximum likelihood approach: $p(\{T\}|\{\}) = \prod_{i=1}^{N} p_i(Y_i | X_i) = \prod_{i=1}^{N} N(Y_i | X_i \beta, \sigma_i^2)$

  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ↳ prob. depends on instance
  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ↳ instances are independent

$$-\log p(\{T\}\{\}) = \sum_{i=1}^{N} \frac{1}{2} \frac{(Y_i - X_i \beta)^2}{\sigma_i^2} + \sum_{i=1}^{N} \frac{1}{2} \log \sigma_i^2 + \underbrace{\sum_{i=1}^{N} \log \sqrt{2\pi}}_{\text{independent of model} \Rightarrow \text{drop}}$$

- 2 cases: $\sigma_i^2$ is known or not know
- case 1: - $\sigma_i^2$ is known, e.g. from a pilot experiment that characterises the measurement device

  $\Rightarrow \sum_{i=1}^{N} \frac{1}{2} \log \sigma_i^2$ is constant, i.e. independent of the model $\beta \Rightarrow$ drop

define the covariance matrix as diagonal matrix

weighted least squares objective:

$$\boxed{\hat{\beta} = \arg\min_{\beta} \sum_{i=1}^{N} \frac{(Y_i - X_i \beta)^2}{\sigma_i^2} = \underset{\arg\min}{(Y - X\beta)^T \Sigma^{-1} (Y - X\beta)}}$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \mathbb{0} \\ & & \ddots & \\ \mathbb{0} & & & \sigma_N^2 \end{pmatrix}$$

$$\frac{\partial \text{loss}}{\partial \beta} = -2 X^T \Sigma^{-1}(Y - X\beta) \overset{!}{=} 0 \Rightarrow \text{weighted normal equations} \quad \boxed{(X^T \Sigma^{-1} X)\beta = X^T \Sigma^{-1} Y}$$

formal solution : weighted pseudo-inverse

$$\hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$$

weighted pseudo-inverse

reduce to OLS by scaling : $\quad \tilde{X}_i = \dfrac{X_i}{\sigma_i} \qquad \tilde{Y}_i = \dfrac{Y_i}{\sigma_i}$

$$\hat{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y} \qquad \text{OLS of new variables}$$

$\Rightarrow$ data preparation: scaling & centering $\quad \tilde{X}_i = \dfrac{X_i - \bar{X}}{\sigma_i} \qquad \tilde{Y}_i = \dfrac{Y_i - \bar{Y}}{\sigma_i}$

weighted pseudo-inverse also works if $\Sigma$ is not diagonal, i.e. the
noise between different instance can be correlated, but we must know correlation

· case 2 $\quad \sigma_i^2$ are unknown $\Rightarrow$ we must learn them along with $\beta$
  more model parameters $\Theta = \{ \beta, \{\sigma_i^2\}_{i=1}^N \}$
  mixed supervised learn and unsupervised learning: we have
  supervision for $Y_i \approx X_i \beta$ , but no training information for $\sigma_i^2$

· $\sum\limits_{i=1}^N \dfrac{1}{2} \log \sigma_i^2$ is now dependent on the model and cannot be dropped

$\Rightarrow$ heteroscedastic loss , Dawid-Sebastiani score

$$\hat{\beta}, \{\hat{\sigma}_i^2\}_{i=1}^N = \underset{\beta, \{\sigma_i^2\}}{\arg\min} \; \sum_i \left[ \dfrac{(Y_i - X_i \beta)^2}{\sigma_i^2} + \tfrac{1}{2} \log \sigma_i^2 \right]$$

Solve by <u>alternating optimization</u> : fundamental approach when the optimization
cannot be solved analytically, but the unknow parameters can be split into
two groups ( here : $(\beta, \{\sigma_i^2\}) = (\theta_1, \theta_2)$ )

generic strategy : ⓪ initialization : make an initial guess $\theta_2^{(0)}$

① for $t = 1, \dots, T$

    I. optimize for $\theta_1^{(t)}$, treating $\theta_2^{(t-1)}$ as constants

    II. optimize for $\theta_2^{(t)}$, treating $\theta_1^{(t)}$ as constants

advantage: I and II are often much simpler than the full optimization,
       may be analytically tractable

disadvantage: the final solution is generally not the global optimum of full problem

application to hetoroscedastic loss : <u>Iteratively Reweighted Least Squares</u>

abbreviate : $\sigma_i^2 = \tau_i$    ⓪ initialization : $\tau_i^{(0)} = 1$

    ① for $t = 1, \dots, T$ :

       I : solve for $\beta$, keeping $\tau_i^{(t-1)}$ fixed

$$\Rightarrow \text{ weighted LSQ } \quad \beta^{(t)} = \arg\min_{\beta} \sum_{i=1}^{N} \frac{(y_i - x_i \beta)^2}{\tau_i^{(t-1)}}$$
(case 1)

       II solve for $\tau_i$, keeping $\beta^{(t)}$ fixed

       define residual : $v_i = y_i - x_i \beta^{(t)}$

$$\{\tau_i^{(t)}\} = \arg\min_{\{\tau_i\}} \sum_{i=1}^{N} \left[ \frac{v_i^2}{\tau_i} + \log \tau_i \right]$$

$$\frac{\partial Cost}{\partial \tau_i} = -\frac{v_i^2}{\tau_i^2} + \frac{1}{\tau_i} \overset{!}{=} 0 \quad\Rightarrow\quad \hat{\tau}_i^{(\epsilon)} = v_i^2$$

- how to choose $T$ ($\#$ of iterations)? - IRCS in theory requires only $\bar{T} = 2$, but due to possible numerical inaccuracies, bit bigger $T$ is safer $T = 5, T = 10$
  - in general for alternating optimization: numerical analysis may provide a theory about convergence speed $\Rightarrow$ $\#$ choose $T$ accordingly / otherwise trial & error

- heteroscedastic loss can also be used for non-linear optimization

$$Y_i \approx f(X_i) \pm \epsilon_i \qquad \epsilon_i \sim N(0, \sigma_i^2 = g(X_i))$$

$\quad\quad\quad \hookleftarrow$ train two non-linear predictors ___↑   e.g. neural networks