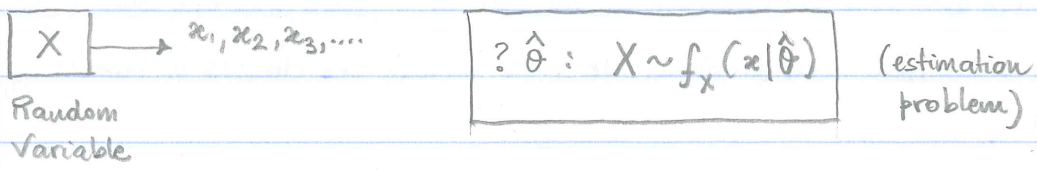


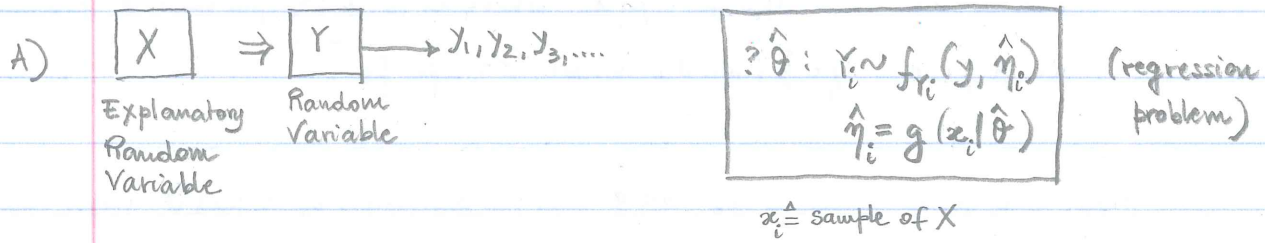
LECTURE 4

So far, we have considered the following problem:



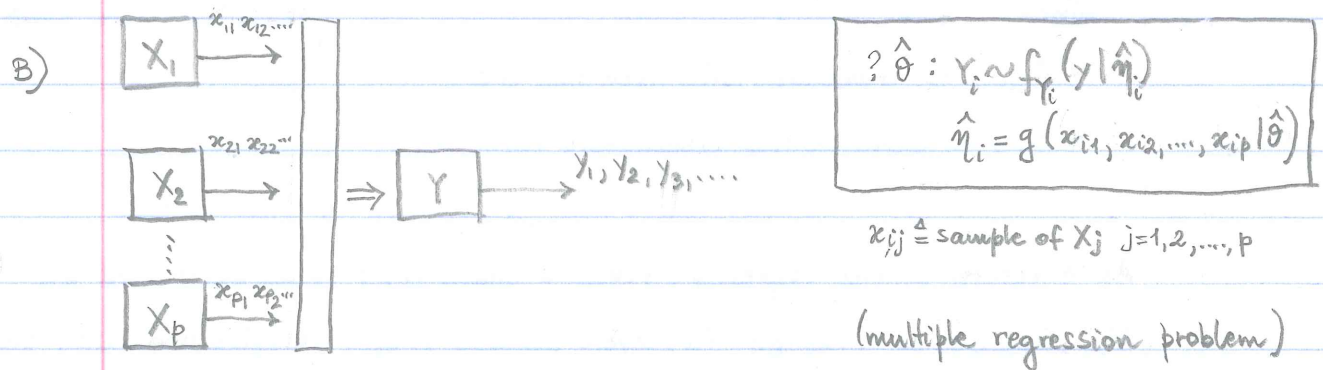
We have approached this problem by using ML estimation, i.e., $\hat{\theta} = \arg \max_{\theta} l(\theta)$
 with: $l(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n \log(f_X(x_i | \theta))$

Let us now consider a more complicated version of this problem:



In this problem, we assume that the vector of parameters η_i is a function of the value of RV $X \Rightarrow$ We must estimate the parameter vector $\hat{\theta}$ that uniquely defines the function between x_i and $\eta_i \forall i$.

The problem can be also posed, as a multivariable problem:



(2)

Both in problem A) and B), functions $f_{Y_i}(\cdot)$ and $g(\cdot)$ can be linear or nonlinear, and determine the existence of one solution $\hat{\theta}$, more than one solution, or no solution.
 \Rightarrow To solve these problems, we must provide the classes of models $f_{Y_i}(\cdot)$ and $g(\cdot)$

EX: Let us assume: $Y_i \sim N(\mu_i, \sigma^2)$ with μ_i - unknown
 σ^2 - known \Rightarrow We set: $\eta_i \triangleq \mu_i$

$$\Rightarrow f_{Y_i}(y | \eta_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\eta_i)^2}{2\sigma^2}}$$

Let us assume: $\eta_i \triangleq \beta_0 + \beta_1 x_i$, where x_i is a sample of the explanatory RV X , i.e.: $\eta_i = g(x_i | \theta)$, where $\theta = [\beta_0 \beta_1]^T$.

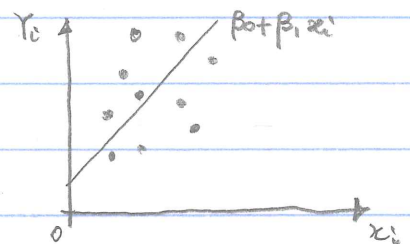
Note: The formulation of f_{Y_i} and g , and the definition of η_i and θ imply that we have:

$$Y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\mu_i} + \epsilon_i$$

RV of \uparrow the i -th evaluation \uparrow RV $\sim N(0, \sigma^2)$

\Rightarrow Linear regression is a special case of "regression problem" that is obtained when:

- * $f_{Y_i}(\cdot)$ is Gaussian with fixed variance
- * $g(\cdot)$ is linear in the explanatory variables



As a result, Linear regression requires that: (i) Y varies unbounded as the sample x of X changes, (ii) every sample y_i of Y is extracted with the same variance (i.e., $\text{var}(\epsilon_i) = \sigma^2 \forall i$), and (iii) the responses Y_i are all normal

⇒ We want to relax requirements (i)-(iii). In particular, let us assume:

- Y_i is bounded (e.g., discrete variable between 0 and n_i)
 - variance of Y_i changes with $i=1, 2, 3, \dots$
 - Y_i is not Gaussian
- } ⇒ For instance: $Y_i \sim B(n_i, p_i)$

Where: $n_i \triangleq$ number of trials used to estimate Y_i

$p_i \triangleq E\left(\frac{Y_i}{n_i}\right)$ - expected value of count Y_i over n_i trials

In this case, we need to identify " η_i ", " $g(\cdot)$ ", and " θ ". Note this:

$0 \leq p_i \leq 1 \forall i \Rightarrow \frac{p_i}{1-p_i} \in [0, +\infty) \Rightarrow \underbrace{\log\left(\frac{p_i}{1-p_i}\right)}_{\triangleq \text{logit}(p_i)} \in (-\infty, +\infty)$, i.e., Logit could vary unbounded as a function of an explanatory variable x_i

⇒ Hence, consider the following case:

$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_i \Rightarrow p_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \Rightarrow$ We can assume:

$\left. \begin{aligned} \eta_i &\triangleq p_i ; \theta \triangleq [\beta_0 \beta_1]^T \\ g(x_i | \theta) &\triangleq \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \end{aligned} \right\} \Rightarrow \boxed{\begin{aligned} Y_i &\sim f_{Y_i}(y | \eta_i) \\ \eta_i &= g(x_i | \theta) \end{aligned}} \quad \begin{array}{l} \text{regression} \\ \text{problem} \end{array}$

where: $f_{Y_i}(y | \eta_i) \stackrel{\text{definition of } B(n_i, p_i)}{\triangleq} P(Y_i = y) = \binom{n_i}{y} \eta_i^y (1 - \eta_i)^{n_i - y}$

Interestingly, the problem can also be formulated in the compact form:

$Y_i \sim f_{Y_i}(y | \eta_i) \quad (*)$
 $\text{logit}(\eta_i) = [1 \ x_i] \theta$

④

Similarly one can formulate the correspondent multiple regression problem:

$$Y_i \sim f_{Y_i}(y | \eta_i) \quad (**)$$
$$\text{logit}(\eta_i) = [1 \ x_{i1} \ \dots \ x_{ip}] \theta$$

Problems (*) and (**) are called "logistic regression" problems.

Note: In both problems, $f_{Y_i}(\cdot)$ is non-Gaussian and $g(\cdot)$ is nonlinear. However $g(x|\eta)$ can be put in a form that is LINEAR in the parameters. Moreover, both problems can be solved by using the ML method:

$Y_i \sim B(n_i, \eta_i) \quad i=1, 2, 3, \dots, n \Rightarrow$ log-likelihood function:

$$l = l(\eta_1, \eta_2, \dots, \eta_n) = \sum_{i=1}^n \log \left(\eta_i^{y_i} (1 - \eta_i)^{n_i - y_i} \right)$$
$$= \sum_{i=1}^n \left[y_i \log \eta_i + (n_i - y_i) \log (1 - \eta_i) \right]$$

By replacing $\eta_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$, we have:

$$l = l(\theta) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 x_i) - y_i \log (1 + e^{\beta_0 + \beta_1 x_i}) + \right. \\ \left. - (n_i - y_i) \log (1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

$\theta = [\beta_0 \ \beta_1]^T$

$$l(\theta) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 x_i) - n_i \log (1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

\Rightarrow The log-likelihood function is concave \Rightarrow The MLE of the parameter vector θ is obtained by imposing $l'(\theta) = 0$ and the S.E. of the ML estimation $\hat{\theta}$ can be obtained from $1/l''(\theta)$

* What are the advantages of Logistic regression vs. Linear regression?

Logistic	Linear
$Y_i \sim B(n_i, \eta_i)$	$Y_i \sim N(\eta_i, \sigma^2)$
$\text{logit}(\eta_i) = \beta_0 + \beta_1 x_i$	$\eta_i = \beta_0 + \beta_1 x_i$

- The pdf is NOT Gaussian (\Rightarrow good to model discrete/quantized events) and the variance needs NOT to be constant across trials
- The function $g(\cdot)$ needs NOT to be linear \Rightarrow It can be part of the more general class of nonlinear functions (e.g., exponential, logarithmic, etc.)
- Despite the nonlinearities, the logistic regression problem can be solved by using the ML method and the solution is unique (i.e., the log-likelihood function is concave)

Note: The advantages reported above are NOT exclusive to the pair $(B(\cdot, \cdot), \text{logit}(\cdot))$
 \Rightarrow We can consider numerous alternatives and still preserve the advantages \Rightarrow
 We can define new types of regression problems. For instance:

(***)

$$Y_i \sim B(n_i, \eta_i)$$

$$\phi^{-1}(\eta_i) = \beta_0 + \beta_1 x_i$$

\uparrow
 $\triangleq \text{probit}(\cdot)$

where: $\Phi(z) \triangleq P(Z \leq z)$ and $Z \sim N(0, 1)$

$$\text{i.e., } \Phi(z) \triangleq \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$$

Note: This choice is not casual but it rather reflects a property of the logistic regression:

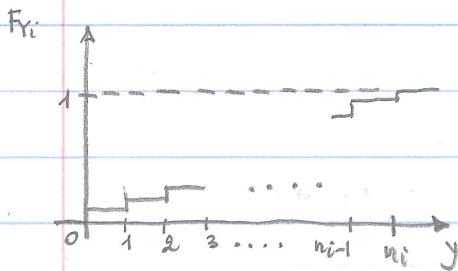
$$p_i \triangleq E\left(\frac{Y_i}{n_i}\right) \Rightarrow F_{Y_i}(y) = \sum_k \binom{n_i}{k} p_i^k (1-p_i)^{n_i-k} \Rightarrow F_{Y_i} \text{ is a combination of}$$

\uparrow CDF

\uparrow number of integers $< y$

expt functions

⑥



where $\text{expit}(p) \triangleq \frac{e^p}{1+e^p} \Rightarrow F_{Y_i}(\cdot)$ has

a sigmoidal shape \Rightarrow Hence, model (***) simply replaces a sigmoidal cdf with another sigmoidal cdf \Rightarrow This may help to better fit experimental data

Ex.: (Latent Variable) Let us assume that the value of a binary RV Y depends on the value of another (continuous) RV W (e.g., Y may model the "perception" event while W captures the intensity that is required to the perception process in order to generate an event):

$$Y = \begin{cases} 1 & W > c \\ 0 & W \leq c \end{cases} \quad \text{where } c \text{ - given and } W \sim N(\mu_W, 1)$$

By defining $Z \triangleq \mu_W - W$ we have: $P(Y=1) = P(W > c) = P(Z < \mu_W - c)$ and $Z \sim N(0, 1)$. Hence, we can describe the process with the model:

$$\begin{cases} Y \sim B(1, \eta) \\ \Phi^{-1}(\eta) = \beta_0 + \beta_1 Z \end{cases}$$

Similarly, one can replace the probability function $B(n_i, \eta_i)$ with another probability function and still preserve the same advantages. For instance, consider:

$$Y_i \sim \mathcal{P}(\eta_i)$$

$$\log \eta_i = [1 \ x_i] \theta$$

$$\text{where: } \theta \triangleq [\beta_0 \ \beta_1]^T$$

$$\eta_i \triangleq \lambda_i$$

$$f_{Y_i}(y|\eta) \triangleq \frac{e^{-\eta} \eta^y}{y!}$$

In this case, by applying the ML method, we have:

$$l(\eta_1, \eta_2, \dots, \eta_n) \triangleq \sum_{i=1}^n \log \left(e^{-\eta_i} \eta_i^{y_i} \right) = -\sum_{i=1}^n \eta_i + \sum_{i=1}^n y_i \log \eta_i$$

$$\Rightarrow l(\theta) = -\sum_{i=1}^n e^{\beta_0 + \beta_1 x_i} + \sum_{i=1}^n y_i (\beta_0 + \beta_1 x_i)$$

$$\uparrow$$

$$\eta_i = e^{\beta_0 + \beta_1 x_i}$$

$$\triangleq -e^{\beta_0} \sum_{i=1}^n e^{\beta_1 x_i} + n\beta_0 \bar{y} + \beta_1 \sum_{i=1}^n x_i y_i$$

$\bar{y} \triangleq$ sample mean

$$\text{Hence, we have: } \frac{\partial l}{\partial \beta_0} = 0 \Leftrightarrow \left(-\sum_{i=1}^n e^{\beta_1 x_i} \right) e^{\beta_0} + n\bar{y} = 0 \Leftrightarrow \hat{\beta}_0 = \log \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n e^{\beta_1 x_i}}$$

$$\frac{\partial l}{\partial \beta_1} = 0 \Leftrightarrow \left(-\sum_{i=1}^n x_i e^{\beta_1 x_i} \right) e^{\beta_0} + \sum_{i=1}^n x_i y_i = 0$$

$$\Leftrightarrow \frac{\sum_{i=1}^n x_i e^{\beta_1 x_i}}{\sum_{i=1}^n e^{\beta_1 x_i}} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n y_i} \quad \text{- A solution can be found numerically} \quad \square$$

Note this: In every example considered thus far, the possibility of using the ML method stems from the fact that (i) the function $f_{Y_i}(\cdot)$ is a combination of exponential functions and (ii) the function $g(\cdot)$ is manipulated to get a new relationship that is linear in the parameters. Finally, note that in all these examples, the variable η_i is the expected value of the data \Rightarrow We can generalize:

8

$$(a) \begin{cases} Y_i \sim f_{Y_i}(y|\eta_i) = h(y) e^{(A(\eta_i)T(y) - B(\eta_i))} \\ g^{-1}(\eta_i) = \underbrace{[1 \ x_i]}_{\text{linear predictor}} \theta \end{cases}$$

Model (a) is called "Generalized Linear Model" (GLM) and function $g(\cdot)$ is called the "link" function, as it links the random and systematic components of the model. In fact:

$$\begin{aligned} \bullet Y_i \sim \mathcal{P}(\eta_i) &\Leftrightarrow f_{Y_i}(y|\eta_i) = \frac{1}{y!} \eta_i^y e^{-\eta_i} = h(y) e^{\underbrace{y \log \eta_i}_{T(y)} - \underbrace{\eta_i}_{B(\eta_i)}} \\ \text{and} & \\ g^{-1}(\cdot) = \log(\cdot) &\Leftrightarrow g(\cdot) = \exp(\cdot) \quad h(y) \end{aligned}$$

$$\begin{aligned} \bullet Y_i \sim \mathcal{B}(n_i, \eta_i) &\Leftrightarrow f_{Y_i}(y|\eta_i) = \binom{n_i}{y} \eta_i^y (1-\eta_i)^{n_i-y} = \\ \text{and} & \\ g^{-1}(\cdot) = \text{logit}(\cdot) & \\ &= \binom{n_i}{y} e^{y \log \eta_i + (n_i-y) \log(1-\eta_i)} \\ &= \binom{n_i}{y} e^{\underbrace{y(\log \eta_i - \log(1-\eta_i))}_{A(\eta_i)} + \underbrace{n_i \log(1-\eta_i)}_{B(\eta_i)}} \\ &\quad \uparrow \quad \uparrow \\ &\quad h(y) \quad T(y) \end{aligned}$$

$$\begin{aligned} \bullet Y_i \sim \mathcal{N}(\eta_i, \sigma^2) &\Leftrightarrow f_{Y_i}(y|\eta_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\eta_i)^2}{2\sigma^2}} \\ &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{y(y-2\eta_i)}{2\sigma^2} + \frac{\eta_i^2}{2\sigma^2}} \\ &\quad \uparrow \quad \uparrow \quad \uparrow \\ &\quad h(y) \quad T(y) \quad A(\eta_i) \quad B(\eta_i) \end{aligned}$$

In this way, regression methods are extended to exponential families (e.g., binomial, Poisson, normal, inverse Gaussian distributions, etc.). Moreover, the choice of the link function is NOT dictated by the chosen exponential family

The introduction of the GLMs leads to the following question: for a given pair (GLM, link function), how to choose the size of θ ?

In this, one can consider the likelihood function $L(\theta)$ and exploit the fact that the higher $L(\theta)$, the better the fit \Rightarrow One can choose a set of parameters θ_1 over another set θ_2 based on the ratio:

$$\frac{L(\theta_1)}{L(\theta_2)} \Leftrightarrow \text{or, equivalently, } r = \log(L(\theta_1)) - \log(L(\theta_2))$$

In particular, if $\theta_1 = \beta_0$ and $\theta_2 = [\beta_0 \beta_1]^T$ we have:

$$D_n \triangleq -2 \log L(\beta_0) \quad \text{- null deviance}$$

$$D_r \triangleq -2 \log L(\theta_2) \quad \text{- residual deviance}$$

$$d \triangleq -2 \log(r) = D_n - D_r \Rightarrow \text{The decision about whether or not } \beta_1 = 0 \text{ depends on the value of the log-likelihood ratio}$$

↑
log-likelihood ratio

One last generalization from the linear regression case:

- In linear regression, we assume that the variance remains the same $\forall i$ $\left\{ \begin{array}{l} \Rightarrow \text{The sample estimation is: } s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ \text{with: } \hat{y}_i = \hat{\eta}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \end{array} \right.$

\Rightarrow The sum: $SSE \triangleq \sum_{i=1}^n (y_i - \hat{y}_i)^2$ (sum of squares for error) accounts for the prediction error

Analogously, one can define:

$$SST \triangleq \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{(total sum of squares)} \Rightarrow \text{It accounts for the entire variability in the data}$$

with $\bar{y} \triangleq$ sample mean

(10)

The ratio: $R^2 \triangleq \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$ is the proportion of variability in Y that is attributable to the regression line

$\Rightarrow R^2$ is the proportion of variability in Y that is explained by X

The definition of R^2 , however, is not valid in the GLM case \Rightarrow An alternative is the Nagelkerke R^2 measure:

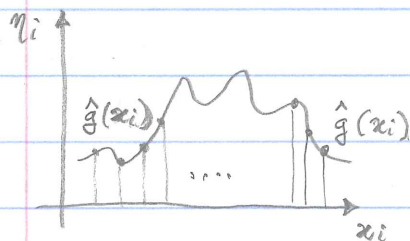
$$R_N^2 \triangleq 1 - \left(\frac{L(\beta_0)}{L(g)} \right)^{2/n}$$

with the max value given for: $R_N^2 \max = 1 - (L(\beta_0))^{2/n}$

* Nonparametric Regression

Consider the original model we started from: $Y_i \sim f_{Y_i}(y|\eta_i)$
 $\eta_i = g(x_{1i}, x_{2i}, \dots, x_{pi})$

In this case, let us assume that the class of functions for $g(\cdot)$ is NOT defined, i.e., there is no set of parameters θ to estimate and a general form for $g(\cdot)$ must be determined \Rightarrow It is a "nonparametric regression" problem



Option #1: We can define $g(\cdot)$ as the sequence of values obtained for each eval. of the variables x 's. For instance:

$g(x)$ is defined by $(\hat{g}(x_1), \dots, \hat{g}(x_n))$
where $\hat{g}(\cdot)$ is an estimation

$(\hat{g}(x_1), \hat{g}(x_2), \dots, \hat{g}(x_n))$ is a "smoother"

One possible smoother is obtained by assuming that each estimation $\hat{g}(x_i)$ is a linear combination of all the measurements y_1, y_2, \dots, y_n according to some weights $(h_{i1}, h_{i2}, \dots, h_{in})$, i.e.,

$$\hat{g}(x_i) = [h_{i1} \ h_{i2} \ \dots \ h_{in}] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Hence, in matrixial form:

$$[\hat{g}(x_1) \ \hat{g}(x_2) \ \dots \ \hat{g}(x_n)]^T = Hy$$

$$y \triangleq [y_1 \ y_2 \ \dots \ y_n]^T$$

$$H = \{h_{ij}\} - n \times n \text{ matrix}$$

Note this: In linear regression we have $\eta_i = [1 \ x_i] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \Rightarrow$ Denoted with i

$$\left. \begin{array}{l} X \triangleq \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \eta \triangleq \begin{bmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_n \end{bmatrix} \\ \theta \triangleq [\beta_0 \ \beta_1]^T \end{array} \right\} \text{ We have: } \eta = X\theta$$

An estimator $\hat{\theta}$ of θ is computed by using the least-squares method and is given

$$\text{by: } \hat{\theta} = (X^T X)^{-1} X^T y \Rightarrow \text{Hence we have: } \hat{\eta} \triangleq [\hat{g}(x_1) \ \dots \ \hat{g}(x_n)] = \underbrace{X (X^T X)^{-1} X^T}_H y$$

H - "hat" matrix

Moreover, since $\eta_1, \eta_2, \dots, \eta_n$ are predictions, one can also compute the variance of the estimation:

$$\text{var}([\hat{g}(x_1) \ \hat{g}(x_2) \ \dots \ \hat{g}(x_n)]^T) = \text{var}(Hy) = H \text{var}(y) H^T = \sigma^2 H H^T$$

↑
If the variance is the same $\forall i=1, 2, \dots, n$

12

Option #2: We can approximate $g(x)$ with a polynomial function of adequate order s :

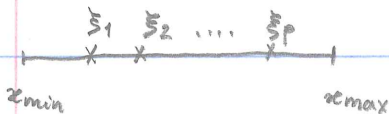
$$g(x) \cong b_0 + b_1 x + b_2 x^2 + \dots + b_s x^s \quad (b)$$

In this case, we can define $w_1 \triangleq \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $w_2 \triangleq \begin{bmatrix} x_1^2 \\ x_2^2 \\ \vdots \\ x_n^2 \end{bmatrix}$, ..., $w_s \triangleq \begin{bmatrix} x_1^s \\ x_2^s \\ \vdots \\ x_n^s \end{bmatrix}$ and have:

$$[\hat{g}(x_1) \quad \hat{g}(x_2) \quad \dots \quad \hat{g}(x_n)] = [1 \quad w_1 \quad w_2 \quad \dots \quad w_s] \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_s \end{bmatrix}$$

w_1, w_2, \dots, w_s are called "basis functions" \Rightarrow We can use the least-squares method to estimate $[b_0 \quad b_1 \quad \dots \quad b_s]^T$ from data $[y_1 \quad y_2 \quad \dots \quad y_n]^T$ as done above.

A variation of the option with basis functions is given by "splines": The interval $x \in [x_{\min}, x_{\max}]$ is divided into $p+2$ sub-intervals (no overlap) and formula (b) is used in each interval with different values $[b_0 \quad b_1 \quad \dots \quad b_s]^T$ and a continuity constrain for the extremes of adjacent intervals, i.e.:



$$(x - \xi_i)_+ \triangleq \begin{cases} x - \xi_i & \text{if } x \geq \xi_i \\ 0 & \text{otherwise} \end{cases}$$

$$\forall x \in [x_{\min}, x_{\max}] \quad i = 1, 2, \dots, p$$

$$g(x) \cong [b_0 \quad b_1 \quad \dots \quad b_s] \begin{bmatrix} 1 \\ x \\ \vdots \\ x^s \end{bmatrix} + \sum_{i=1}^p [b_{s+(i-1)s+1} \quad \dots \quad b_{2s+(i-1)s}] \begin{bmatrix} (x - \xi_i)_+ \\ \vdots \\ (x - \xi_i)_+^s \end{bmatrix}$$

The function $g(x)$ is polynomial on each segment $[\xi_i, \xi_{i+1}]$ and, defined

the variables: $x_1 \triangleq x; x_2 \triangleq x^2; \dots; x_s \triangleq x^s; x_{s+1} \triangleq (x - \xi_1)_+; \dots$
 $x_{2s} \triangleq (x - \xi_1)_+^s; \dots; x_{ps+1} \triangleq (x - \xi_p)_+; \dots; x_{(p+1)s} \triangleq (x - \xi_p)_+^s,$

one can estimate the parameters $b_0, b_1, b_2, \dots, b_{(p+1)s}$ by using a multiple linear regression method, under the assumption that the variance, σ^2 is the same at each sample and the samples are independent.

Note: The solution with basis functions (in particular, splines) can be easily extended to the link function of a GLM, i.e., we can consider:

$$Y_i \sim f_{Y_i}(y | \eta_i)$$

$$g^{-1}(\eta_i) = W \theta \quad \text{where } W \text{ is a matrix of basis functions:}$$

$$W \triangleq \begin{bmatrix} 1 & w_{11} & w_{12} & \dots & w_{1s} \\ 1 & w_{21} & w_{22} & \dots & w_{2s} \\ \vdots & & & \ddots & \\ 1 & & & & w_{ss} \end{bmatrix} \quad w_{ij} \triangleq x_i^j$$

Option #3: In the GLM template, η_i is the expected value of the RV $Y_i \Rightarrow$ We can

$$\left. \begin{array}{l} \text{consider: } \eta_i = E(Y_i) \\ \eta_i = g(x_i) \end{array} \right\} \Rightarrow \eta_i = E(Y_i | x_i)$$

This suggests that, if the values x_1, x_2, \dots, x_n are close to a certain value x and $g(\cdot)$ is a smooth function, then it is reasonable to expect that:

$\bar{y} \triangleq \frac{1}{n} \sum_{i=1}^n y_i$ is an estimator for $E(Y_i | X=x) \Rightarrow$ We can consider a neighborhood of x and average the values y_i that are obtained for x_i in that neighborhood \Rightarrow $g(\cdot)$ can be estimated by using the local averages \bar{y} . Also, in a more general form, one can use all y_i to compute \bar{y} , but different weights are given to the

(14)

values y_i depending on the distance of x_i from x (i.e., the farther x_i , the lower the weight w_i):

$$\bar{y}_w = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \quad (c)$$

Formula (c) is particularly useful when $w_i = k\left(\frac{x-x_i}{h}\right)$ $i=1, 2, \dots, n$, where $k(u)$ is a smooth function (e.g., a pdf) and is called **KERNEL**, and h controls the sensitivity to $x-x_i$ and is called **BANDWIDTH** ($k(u) \sim N(0,1) \Rightarrow h$ plays the role of s.d.). In this case, we have:

$$g(x) = \frac{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) y_i}{\sum_{i=1}^n k\left(\frac{x-x_i}{h}\right)}$$

References:

Textbook: ch. 14

ch. 12 (section 12.1, 12.5.1, 12.5.3 up to page 343)

ch. 15 (section 15.1, 15.2.1, 15.2.2, 15.2.3, 15.2.4, 15.2.7, 15.3.1)

Recommended Reading:

McCullagh P. & Nelder J.A. "Generalized Linear Models", 2nd ed., CRC, 1990: read ch. 2 (a copy is on Husky CT)