model: $Y = X\underline{\beta} + \varepsilon$ → **Mean structure** ← true model

- <u>idea</u>: <u>data</u> are <u>generated</u> from an <u>underlying system</u>, which is assumed to have the <u>form</u>: $y = f(x_1, ..., x_m) + \varepsilon$, where $f$ is **<u>unknown</u>**. ← use <u>data to gain information</u>

- <u>regression</u> *approximates* the <u>mean structure</u> $f$ by a <u>linear combination</u> of (<u>known</u>) *base functions* $\underline{g_i}(x_1, ..., x_m)$'s, $i = 1, ..., p$, i.e.,

  $$\underline{f} \longleftarrow \sum_{i=1}^{p} \beta_i \cdot g_i(x_1, ..., x_m)$$
  data determine → unknown
  
  [especially on a "<u>local</u>" region of the predictors]

  ➢ when the structure of $f$ is <u>simple</u> and <u>almost linear</u>, it can be approximated by a <u>simple structure</u> with <u>fewer terms</u>, e.g.,
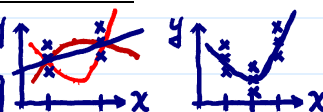
  [complex model] ← 是否有需要?

  $$E(y) = \underline{f} \approx \beta_0 + \underline{\beta_1}\,\underline{x_1} + \cdots + \beta_m\,\underline{x_m}$$

  - 是否有能力? ⊙ **Q**: <u>nature</u> is <u>simple</u>? → *lack of fit problem* → [change $g_i$'s (transformation) / add more $g_i$'s]
  - ⊙ **Q**: are there <u>sufficient data</u> to support/fit a <u>complex model</u>?

  ➢ when $f$ is <u>complex</u> and <u>non-linear</u> ⇒ need <u>more terms</u> to get a <u>good approximation</u> [on a <u>wider</u> region of the predictors]

  [if YES & YES]

  

  - <u>more parameters</u>, need <u>more degrees of freedom</u>, i.e., <u>more data</u>
  - e.g., <u>2 levels</u>, only <u>linear effects</u>; <u>3 levels</u>, <u>linear</u> and <u>quadratic effects</u>
  - ⊙ **Q**: what other <u>complex models</u>? → *What base functions should we consider?* [# of distinct $x_i$'s]

- <u>base functions</u> for <u>quantitative</u> and <u>qualitative</u> predictors $x_i$'s are defined in <u>different ways</u>
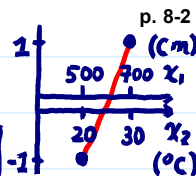
---

[e.g. $x → x - \bar{x}\mathbb{1}$] ← **location and/or scale change**



- $\underline{x_i} → (x_i + a)/b$ or $y → (y + a)/b$, where $a$ and $b$ are given constants.

  $a$: <u>change of location</u>, $b$: <u>change of scale</u>   [**Q**: Can we compare the magnitudes of $\hat{\beta}_i$'s to identify <u>important effects</u>? Ans. In general, <u>NO</u>.]

- **Q**: why we might want to do this?

  $\hat{\beta}$ ⇑ $(x'x)^{-1}$

  ➢ <u>predictors</u> of similar magnitude are easier to compare $\hat{\beta}_i$'s ← [• ridge regression / • LASSO]

  ➢ [or change of units] <u>rescaling</u> may make $\hat{\beta}$ easier to read and may aid interpretability

  $|x'x| \approx 0$

  ➢ <u>numerical stability</u> is <u>enhanced</u> when all predictors are on a <u>similar scale</u> [500 cm  700 cm]

  ➢ for <u>experimental data</u>, it's often that we <u>code two levels</u> (say, $20°C, 30°C$) → $(-1, 1)$; <u>three levels</u> (say, $20°C, 30°C, 40°C$) → $(-1, 0, 1)$ → remove units

- <u>influence</u> caused by <u>location/scale change</u> on $x_i$ (i.e., $x_i → (x_i+a)/b$) $= x'$ $= \frac{1}{b}x + \frac{a}{b}\mathbb{1}$

  ➢ (under a <u>model</u> with <u>intercept</u>) <u>overall</u> $F$-test, $t$-test, $R^2$, $\hat{\sigma}$ all unchanged    [∵same $\Omega$ / ∵same $\omega$]

  $$E(y) = \beta_0 + \cdots + \beta_i x_i = (\beta_0 - \beta_i a) + \cdots + (b\beta_i)\left(\frac{x_i + a}{b}\right) = \beta_0' + \cdots + \beta_i' x_i'$$

  $$t_i = \frac{\hat{\beta}_i}{se(\hat{\beta}_i)} = \frac{b\hat{\beta}_i}{se(b\hat{\beta}_i)} \quad [\hat{\beta}_i']$$

  ➢ $\underline{\hat{\beta}}$ <u>change</u>: $\hat{\beta}_i → b\hat{\beta}_i$ [$\hat{\beta}_i'$], $\hat{\beta}_0 → \hat{\beta}_0 - a\hat{\beta}_i$ ← $\hat{\beta}_0'$

- <u>influence</u> caused by <u>location/scale change</u> on $y$ (i.e., $y → (y+a)/b$) $= y'$ $= \frac{1}{b}y + \frac{a}{b}\mathbb{1}$

  [TSS / RSS]

  ➢ (under a <u>model</u> with <u>intercept</u>) <u>overall</u> $F$-test, $t$-test, $R^2$ unchanged

  $$y = \beta_0 + \sum_i \beta_i g_i(x) + \varepsilon \Rightarrow y' = \frac{y+a}{b} = \frac{\beta_0 + a}{b} + \sum_i \left(\frac{\beta_i}{b}\right)g_i(x) + \frac{\varepsilon}{b} = \beta_0' + \sum_i \beta_i' g_i(x) + \varepsilon' \Leftarrow Var(\varepsilon')$$

  ➢ $\hat{\beta}$ and $\hat{\sigma}$ <u>change</u>: $\underline{\hat{\sigma} → \hat{\sigma}/b}$, $\hat{\beta}_i → \hat{\beta}_i/b$ for <u>each</u> $i$, $\hat{\beta}_0 → (\hat{\beta}_0 + a)/b$ $= \frac{1}{b^2}Var(\varepsilon)$

❖ **Reading**: Faraway (2005, 1st ed.), 5.2  [$\sqrt{Var(\varepsilon)}$  $\hat{\beta}_i'$]   ❖ **Further reading**: D&S, 16.2, 16.3 [$\hat{\beta}_0'$]

base functions are polynomial terms ⟶ **Polynomial regression** ⌐ unknown true $f$
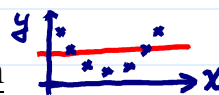
e.g. differentiable infinitely. no jump points no broken lines

- **Q**: when to use polynomial regression? ⇒ the relationship between response and *quantitative* predictors is **smooth**, but not a straight line.

  ⌐ $E(y_x)$     ⌐ $x_1, .., x_m$

- **idea supports the approach** ⇒ any smooth function (mean structure of the underlying system) can be approximated by a polynomial of high enough degree

- one predictor case:    ⌐ locally

  $$y = \beta_0 + \beta_1 \underline{x} + \beta_2 \underline{x^2} + ... + \beta_d \underline{x^d} + \varepsilon$$
  
  ⌐ base functions ⟶

  Taylor's expansion

  ➤ choice of $d$

    - start with $y=\beta_0+\beta_1 x$, keep adding polynomial terms until last term added is not significant. ⇒ danger: stop too soon

    consider it's for
    • prediction
    • interpretation

    - start with a large $d$ and recursively remove insignificant largest term

    - use added variable plot/partial residual plot to gain information about $d$

  ➤ **Q**: Consider the model $y=\beta_0+\overset{×}{\beta_1}x+\overset{∨}{\beta_2}x^2$. what if $\beta_1$ not significant, but $\beta_2$ is significant? should $x$ be removed from the model?

    $E(y)=\beta_0+\beta_2 x^2$ has maximum minimum at $x=0$   ← $x^2$

    - $x$ and $x^2$ could be highly correlated

    - location shift: $x \to x+c \Rightarrow \hat{\beta}_2$ unchanged, but $\hat{\beta}_1$ may become significant

    $E(y)=\beta_0 + \overset{f\approx 0}{\hat{\beta}_1 x} + \beta_2 x^2 = \beta_0' + \beta_1'(x+c) + \beta_2'(x+c)^2$ ← do not want our model sensitive to the change $x \to x+c$

    - recommendation: do not remove insignificant lower-order terms from model when the highest-order term is significant

- two predictors $x_1, x_2$ case: ⌐ 1st-order model ($d=1$)

  $$y = \beta_0 + \beta_1 \underline{x_1} + \beta_2 \underline{x_2} + \beta_{11} \underline{x_1^2} + \beta_{22} \underline{x_2^2} + \beta_{12} \underline{x_1 x_2} + \varepsilon \quad (d=2, \text{2nd-order model})$$

  ➤ the cross-product term $x_1 x_2$ can be interpreted as an "interaction" effect, e.g.,

  cf. $E(y) = \beta_0 + \beta_1 \underline{x_1} + \beta_2 \underline{x_2} + \beta_3 \underline{x_1 x_2}$, where $x_1, x_2 \in \{-1, 1\}$ ← 交互作用

  $x_1 = +1 \Longrightarrow E(y) = \underline{(\beta_0 + \beta_1)} + \underline{(\beta_2 + \beta_3) x_2}$

  $x_1 = -1 \Longrightarrow E(y) = \underline{(\beta_0 - \beta_1)} + \underline{(\beta_2 - \beta_3) x_2}$
  
  difference of slopes $= 2\beta_3$

  # of parameters $= 1 + m + m + \binom{m}{2}$ $= (m+1)(m+2)/2$

  ➤ models for more predictors can be similarly extended

  $$y = \underline{\beta_0 + \sum_{i=1}^m \beta_{1,i} \underline{x_i}} + \sum_{i=1}^m \beta_{2,i} \underline{x_i^2} + \sum_{1 \leq i < j \leq m} \beta_{3,ij} \underline{x_i x_j} + \epsilon$$
  ⌐ 1st-order model

  ➤ increasing degree $d$ ⇒ model may have too many parameters

| # of $x_i$'s | $d=2$ | $d=3$ |
|---|---|---|
| 2 | 6 | 10 |
| 3 | 10 | 20 |
| 4 | 15 | 35 |

- orthogonal polynomials   ⌐∵ collinearity ⇒ $|X^T X| \approx 0 \Rightarrow (X^T X)^{-1}$ unstable

  ➤ polynomial terms can cause numerical instability (especially when $d$ large) and collinearity

  ➤ example: 2nd-order model

  each for different patterns    base functions    change base functions

  true $f$ not a 2nd-order polynomial ⟶ approximate locally using a 2nd-order polynomial

  $= \beta_0$   $+ \beta_1$  ←$x$  $+ \beta_2$  ←$x^2$

  $= \beta_0$   $+ \beta_1$   $a+bx$   $+ \beta_2$   $a'+b'x+c'x^2$

➢ define $z_0=1$, $z_1=a_1+b_1x_1$, $z_{11}=a_2+b_2x_1+c_2x_1^2$, $z_{111}=a_3+b_3x_1+c_3x_1^2+d_3x_1^3$, ...

Find $a_i$, $b_i$, $c_i$, ..., that make $z_j^Tz_k=0$ if $j\neq k$ (and $\|z_i\|=1$ sometimes)

scale change → $\Omega\ni\hat{y}^T\epsilon\,\Omega^\perp$

**Gram-Schmidt process**

$1 \times x^2 x^3$

■ can apply regression to obtain $z_0$, $z_1$, $z_{11}$,... (note: $\hat{y}^T\hat{\varepsilon}=0$), e.g., regress $x_1$ on $z_0$, then the residuals is proportional to $z_1$; regress $x_1^2$ on $z_0$, $z_1$ and the residuals is proportional to $z_{11}$. In R, built-in function is provided to construct orthogonal polynomials.

■ cross-product terms (i.e., interactions) can be defined in a similar manner (e.g., $z_{12}=a+bx_1+cx_2+dx_1x_2$, regress $x_1x_2$ on $z_0$, $z_1$, $z_2$, and the residuals are proportional to $z_{12}$)

$1,2,z_{11},z_{111}$

• $x_1=\hat{\beta}_0 z_0+\hat{\varepsilon}_1$
$\Rightarrow \hat{\varepsilon}_1=x_1-\hat{\beta}_0 z_0$: a 1st-order polynomial of $x_1$
$\Rightarrow \hat{\varepsilon}_1\perp z_0$
$\Rightarrow \hat{\varepsilon}_1\propto z_1$

• $x_1^2=\hat{r}_0 z_0+\hat{r}_1 z_1+\hat{\varepsilon}_2$
$\Rightarrow \hat{\varepsilon}_2=x_1^2-\hat{r}_0 z_0-\hat{r}_1 z_1$: a 2nd-order polynomial of $x_1$
$\Rightarrow \hat{\varepsilon}_2\perp z_0$
　 $\hat{\varepsilon}_2\perp z_1$
$\Rightarrow \hat{\varepsilon}_2\propto z_{11}$

➢ change model based on polynomial terms to model based on $z$'s, e.g.,

**In DOE, often use Z's, rather than X's**

$y=\beta_0+\beta_1 x_1+\beta_2 x_2+\beta_{11} x_1^2+\beta_{22} x_2^2+\beta_{12} x_1x_2+\epsilon$
$\longrightarrow y=\beta_0'+\beta_1' z_1+\beta_2' z_2+\beta_{11}' z_{11}+\beta_{22}' z_{22}+\beta_{12}' z_{12}+\epsilon$

→ factor X. (A,B,C)→(-1,0,1), (A,B,C)→(1,-2,1)

← interested in $\beta'$, not $\beta$

the two models have same column space $\Omega$ (i.e., same $R^2$, $\hat{\sigma}$, overall $F$), but interpretation of $\beta$'s and $\beta'$'s are different (i.e., different estimates, $t$-tests)

➢ orthogonality can save works when selecting model (do not have to refit after deleting term), it's more convenient for fitting and testing ← Why? ∵ orthogonality (LNp. 5-9)

• properties of polynomial model — e.g. curvature & interaction effects

➢ offer more flexible relationship

➢ remember that it's an approximation, we usually do not believe it exactly represents the underlying reality

Taylor expansion on a local area (response surface methodology)

---

➢ polynomials have the advantage of smoothness → polynomial: infinitely differentiable & $\exists k$ s.t. $\frac{d}{dx^t}E(y)=0$, $t\geq k$

cf.

➢ but, have the disadvantage that each data point affects the fit **globally**

solutions of $d/dx\,E(y)=0$

**LNp. 6-6**
**Lab06-03**

➢ For larger values of $d$, the fitted polynomial curves may become wiggly. reason: the curve may capture the random variation, rather than the overall shape of the relationship between predictors and response.
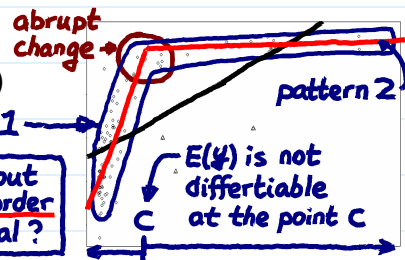
check overfitting (LNp.6-6)

➢ polynomial model is hard to fit "jump function" "local" change

jump

❖ **Reading**: Faraway(2005, 1st ed.), 7.2.2 　❖ **Further reading**: D&S, 12.1, 12.3, 22.2

## broken stick (line) regression (segmented regression)

• Recall. polynomial regression: suitable for smooth mean structure, but cannot capture local abrupt change (example?)

abrupt change →

pattern 2

**Q**: how to relax the smoothness restriction?

pattern 1 →

⇒ one solution: broken line regression.

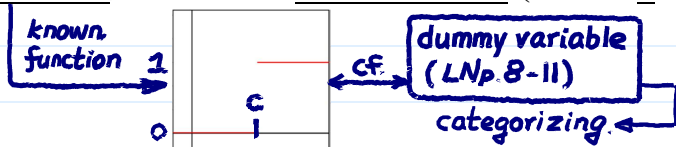How about a 2nd-order polynomial?

$E(y)$ is not differentiable at the point C

• **Q**: when to use broken line regression?

⇒ believe that different regression models apply in different regions of data, and the fit should be continuous at the broken points

→ but not differentiable
⇒ smoothness relaxed

• suppose the break occurs at the *known* value $c$, define the base function (where $c$ is called a *knot*):

$$d_c(x)=\begin{cases} 1, & \text{if } x>c, \\ 0, & \text{if } x\leq c.\end{cases}$$

known function

cf.

dummy variable (LNp. 8-11)

categorizing

p. 8-7

- model: $y = \beta_0 + \beta_1 x + \beta_2 (x-c)d_c(x) + \varepsilon$ — a known function → still a linear model

| intercept and slope of the line on region $x \le c$ | | |
|---|---|---|

$= \beta_0$ ☐ $+ \beta_1$ ☐ $+ \beta_2$ ☐ $+ \varepsilon$

$(c \cdot x)(1 - d_c)$ ← alternative

**Q:** Why can the base function $(x-c)d_c(x)$ reduce "global" influence of data on the fit?

consider its influence on $\hat{\beta} = (X^TX)^{-1}X^TY$

What is the contribution of the term $\beta_2 (x-c)d_c(x)$ on $\hat{y}$?

difference of the 2 slopes $= \beta_2$

$$E(y) = \begin{cases} \beta_0 + \beta_1 x, & \text{if } x \le c, \\ (\beta_0 - \beta_2 c) + (\beta_1 + \beta_2) x, & \text{if } x > c, \end{cases}$$

  ➢ the two lines meet at $c$ ⇒ continuous fit

  ➢ notice only 3 parameters in the model ⇒ one degree of freedom is saved because of the continuity restriction

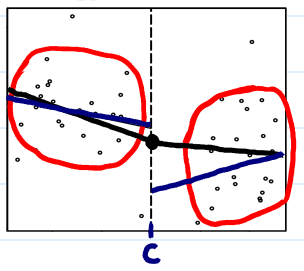⊙ when $c$ is unknown — change point problem

| 4 parameter model |
|---|
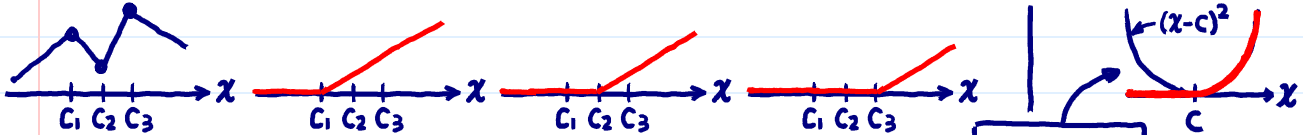→ $E(y) = \beta_0 + \beta_1 d_c(x) + \beta_2 x + \beta_3 (x-c)d_c(x)$  interaction ↲

  ⇒ can regard $c$ as a parameter ⇒ not a linear model any more
  ⇒ can be estimated by nonlinear regression

- generalization: more knot points or more predictors ⇒ define more base functions

$c_1\ c_2\ c_3$  →$x$     $c_1\ c_2\ c_3$ →$x$     $c_1\ c_2\ c_3$ →$x$     $c_1\ c_2\ c_3$ →$x$    $c$ ←$(x-c)^2$

- broken curve regression: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 (x-c)d_c(x) + \beta_4 (x-c)^2 d_c(x) + \varepsilon$

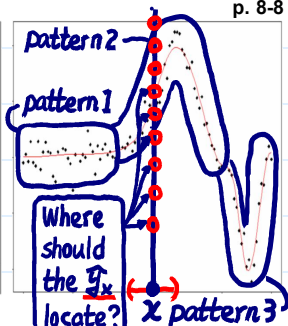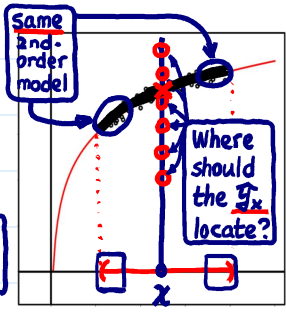❖ **Reading**: Faraway (2005, 1st ed.), 7.2.1    ❖ **Further reading**: D&S, 14.3

# regression spline and LOWESS

p. 8-8

- **concept**: fitting using local points ⇔ better fit

develop $\hat{y}_x$ at any $x$    all data   $y_i$ ← closer → $\hat{y}_i$

  using global points ⇔ smoothness

**Q**: which should you choose for your data?
 using local points or global points?

2nd-order model
$E(y) = \beta_0 + \beta_1 x + \beta_2 x^2$

Same 2nd-order model

Where should the $\hat{y}_x$ locate?

pattern 2
pattern 1
Where should the $\hat{y}_x$ locate?  x pattern 3

- regression spline

  ➢ **concept**: define different *base functions* to fit data
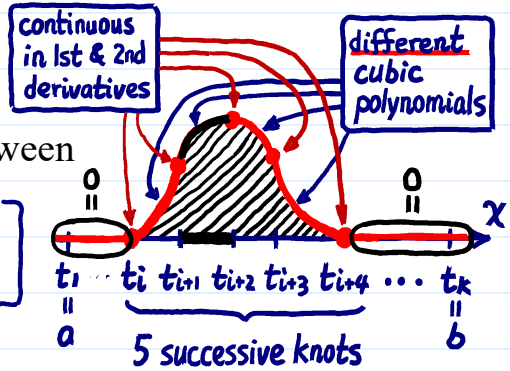
non-differentiable at point C

  ▪ power function: smoothness, but non-zero across the whole range → use global points

Should we include the information in the data of pattern 1 to predict $\hat{y}_x$

  ⊙ broken line: lesser smooth, but localizing the influence of data point

  ▪ B-spline: compromise between smoothness and local influence

  ➢ cubic B-spline base functions: $g_1,..., g_{k-4}$ defined on an interval $[a, b]$ with knot-points $t_1 \le ... \le t_k$ (no need to be equally spaced) satisfying:

    1. non-zero on interval defined by 5 successive knots and zero elsewhere ⇒ local influence

continuous in 1st & 2nd derivatives

different cubic polynomials

    2. a cubic polynomial for each sub-interval between successive knots — infinitely differentiable

    3. continuous, and continuous in its 1st and 2nd derivatives at each knot ⇒ smoothness

$t_i \cdots t_i\ t_{i+1}\ t_{i+2}\ t_{i+3}\ t_{i+4} \cdots t_k$  →$x$
$\parallel$                                              $\parallel$
$a$        5 successive knots         $b$

    4. integrate to one over its support

p. 8-9

**B-spline basis functions**



- base function at the ends of the interval are defined differently to ensure continuity

- regress $y$ on these B-spline base functions, i.e.,

  possesses properties 2 & 3 in LNp.8-8

  $y = \Sigma_i \, \beta_i \, g_i(x) + \varepsilon$     knots are known

  (Note: $g_i$'s are known functions for given $t$'s $\Rightarrow$ it's still a linear model)

non-parametric approach in statistic : dim(parameters) = ∞

non-parametric regression → distribution of $\varepsilon$ / mean structure $E(y_x)$

knot points: 0, 0, 0, 0, 0.2, 0.4, 0.5, 0.6, 0.7, 0.8, 0.85, 0.9, 1, 1, 1, 1

**Q**: why the knots are dense in some region?

- LOWESS (LOcally WEighted Scatterplot Smoothing)

  Note $\dim(\hat{y}) = \dim(\Omega) = \#$ of $\beta_i$'s

  **Spline fit**

  

  - Recall: in previous models, # of parameters are finite

  - nonparametric regressions:

    model: $y = f(x_1, \ldots, x_m) + \varepsilon$

    $f = \sum_{i=1}^{} \beta_i g_i$ , e.g. ↑

    - parametric regression: assume $f$ is from a family of functions, in which # of parameters is finite  cf.

    may consider the mean $E(y_x)$ as a free parameter at every $x$

    - nonparametric regression: assume $f$ is smooth only (# of parameters $= \infty$)

---

p. 8-10

- method: (see example) → at each $x$, $\hat{y}_x = ?$

  
  Step 1

  (i) fix window width
  (ii) fix # of nearest neighbors

  $\hat{y} = ?$  large weight / small weight  $x_0$

  1. select a window ($\Rightarrow$ local information), and a weighting function ($\Rightarrow$ closer points, more contribution)

  Step 2

  weighting function ← weights  $x_0$

  2. use weighted (closer points, higher weights) average of $y_i$'s in the window to estimate fitted value  cf. Recall. WLS

  no sensible estimation of some meaningful $\beta_i$'s

  3. repeat as the window moves

    - width of window is an issue (larger window, smoother curve)

    too large window
    $\Rightarrow$ 規律 → 隨機
    $\Rightarrow$ underfitting

    too small window
    $\Rightarrow$ 隨機 → 規律
    $\Rightarrow$ overfitting

    Step 3

    $y = \beta_0 + \beta_1 x + \varepsilon$
    weighted average $\bar{y}$   $y = \beta_0 + \varepsilon$
    $\hat{\beta}_0 = \bar{y}$

    - can plot fitted value for a variety of widths and pick best result

    - different width of window to estimate $f$ along the range of $x$

    Result
    $\hat{y}_x$   $x_0$

    - sensitive to outliers: use median, not average

  - LOESS: change step 2 to locally weighted (1st or 2nd order) polynomial regression (see example)

  - difficult if extrapolation is required (same difficulty in regression spline) ← cf. → parametric model : $\hat{y}_x = \sum_i \hat{\beta}_i g_i(x)$

  - nonparametric regressions are useful for fitting a curve for residual plots, added variable plots, partial residual plots

❖ **Reading**: Faraway (2005, 1st ed.), 7.2.3; Faraway (2006), chapter 11