

Metoda nejmenších čtverců a její aplikace

(Least square method and its applications)

1 Motivace

Reálné rozdělovací funkce se více či méně odlišují od ideálu normální funkce. Má to více příčin.

- Na rozdíl od normálního rozdělení se v reálných datech objevují mnohem častěji odchýlené body – ty mohou velice silně zkreslit určení polohy středu i rozptyl.
- Zkoumaný vzorek není homogenní – je to směs třeba několika normálních rozdělení (výška vzorku ženy + muži).
- Ve hře je závislost na dalších parametrech, nejčastěji na čase. Např. sinusovitý průběh jasnosti se projeví v bimodálním rozdělení hustoty pravděpodobnosti, algalida má ostrý vrchol v maximu jasnosti a rozsáhlé křídlo apod.

Pokud by měřená veličina měla být v čase konstantní (např. jasnost neproměnné hvězdy), pak by standardní odchylka od průměru měla být dána pouze (!) nepřesností měření dotyčné veličiny (včetně systematických chyb a nestabilit). Tuto nepřesnost lze zhruba odhadnout třeba porovnáním několika po sobě rychle následujících měření. Zjistíme-li, že rozptyl měření určité veličiny je zjevně větší než očekávané nejistoty měření, lze usoudit, že ona veličina zřejmě bude funkcí nějaké další proměnlivé veličiny, nejčastěji času. Proměnné objekty vždy byly a budou předmětem soustředěného zájmu astrofyziků, protože charakterem své proměnnosti toho o sobě prozradí mnohem více, než objekty neproměnné.

Časová závislost měřených veličin (magnetické pole, jasnost, intenzita spektrálních čar, polarizace apod.), hledání trendů, cyklických změn, periodicit apod - to jsou nejčastější úkoly které praktická astrofyzika řeší. Nejoblíbenějším nástrojem pro zpracování těchto závislostí je tzv. *metoda nejmenších čtverců* (MNČ - LSM).

Dříve než přistoupíte k aplikaci MNČ, doporučuji abyste si celou situaci nejprve zevrubně obhlédli, což mj. znamená, že si do nejrůznějších grafů či schémat vynesete vzájemné závislosti všech možných veličin dotyčného objektu, ať už vámi naměřených nebo převzatých z literatury. Věřte, že tyto „obrázky“ vám o povaze vzájemných souvislostí mezi jednotlivými charakteristikami poví více než sebedokonalejší číselné rozbory.

Zjistíte-li, že zobrazené výsledky měření $\{x_i, y_i\}$ jeví jistou závislost, jistě pocítíte neodolatelné nutkání tuto závislost *proložit* (fit) nějakou elegantní hladkou křivkou. Dříve, než se do toho pustíte, byste ale měli zvážit, zda je to skutečně nezbytné! Chceme-li totiž jen doložit, že ona závislost existuje, je poctivější do grafu žádnou křivku nevkruslovat, stačí jen zvolit vhodná měřítka na osách a obrázek prezentovat v jeho originální podobě. Pouze tehdy, chceme-li s výsledky proložení dále pracovat a něco z nich vyvozovat, je případné pustit se do prokládání.

2 Úvod

Jedním z nejčastějších úkolů, s nímž se setkáte při zpracování pozorování, je zjistit a matematicky vyjádřit průběh závislosti jedné pozorované veličiny (y) na jiné pozorované veličině (x).

Velichinu y , kterou zpravidla zjišťujeme s menší relativní přesností (např. hvězdnou velikost), budeme nadále nazývat závislou proměnnou a druhou pozorovanou veličinu, určenou přesněji

(nejčastěji čas), budeme považovat za nezávisle proměnnou veličinu. Závislou a nezávislou veličinu nelze během výpočtu zaměňovat! Reálný vztah mezi oběma veličinami udává (neznámá) funkce $y(x)$.

Dejme tomu, že máme k dispozici celkem n měření dvojic nezávislé a závislé veličiny: $\{x_i, y_i\}$, které mohou být v odůvodněných případech doplněna tzv. *váhou měření* w_i , která kvantifikuje spolehlivost příslušného měření. Zpravidla je tato váha nepřímo úměrná kvadrátu nejistoty určení závislé veličiny y_i .¹

Hledejme nyní takovou funkci $F(x)$, která by co nejlépe odpovídala skutečnému průběhu závislosti $y(x)$, jež je naznačena n dvojicemi $\{x_i, y_i\}$, resp. $\{x_i, y_i, w_i\}$. Triviálním řešením této úlohy je pospojování všech po sobě následujících bodů lomenou čarou, případně nějakou hladkou, dostatečně zvlhčenou čarou (např. polynomem $n-1$ stupně), která by procházela důsledně všemi naměřenými body. Takovýto postup by ovšem přicházel v úvahu snad jen tehdy, kdyby byla poloha jednotlivých bodů grafu známa s absolutně přesně, což je nereálné. Daleko lepší výsledky dává prostá grafická metoda, kdy mezi body vnesenými do grafu táhneme od ruky hladkou křivku, která dle našeho přesvědčení co nejlépe vyjadřuje pozorovanou závislost. Nevýhodou je však to, že tento způsob proložení není obecně reprodukovatelný (i vy sami podruhé proložíte závislost trochu jinak), navíc se s tímto grafickým řešením potom dosti špatně pracuje. Proto dáváme přednost takovým metodám, které vedou k analytickému vyjádření prokládané funkce a k objektivnímu, reprodukovatelnému stanovení kritéria nejlepší shody.

Obvykle postupujeme tak, že si hned na počátku definujeme tzv. *regresní model* (regression model). Regresním modelem si z nekonečného množství funkcí, jimiž by bylo možno pozorovanou závislost proložit, vybereme jen jistou omezenou množinu funkcí, přičemž každá z funkcí této zvolené množiny modelových funkcí bude plně definována g parametry, které si pracovně označíme $\beta_1, \beta_2, \beta_3, \dots, \beta_g$. Veličina g pak vyjadřuje *počet stupňů volnosti* (degree of freedom) zvoleného modelu.

Na tom, zda si umíme již předem vytipovat optimální regresní model, který v sobě obsahuje funkce co nejpodobnější tušené závislosti $y(x)$, závisí úspěch nebo neúspěch celého našeho dalšího počínání. Pokud nevíme o fyzikální podstatě závislosti jedné z pozorovaných veličin na druhé vůbec nic, pak jako regresní model volíme soubor těch co nejjednodušších funkcí - polynomy, harmonické funkce - s nimiž je radost pracovat. Jestliže však již předem víme, jakým typem funkce by měla být pozorovaná závislost popsána, měli bychom to brát ohled, jinak si způsobíme zbytečné problémy při interpretaci zjištěné závislosti.

Volba odpovídajícího regresního modelu s optimálním stupněm volnosti je tím nejdůležitějším momentem při zpracování, momentem na němž ve značné míře závisí i výsledky a jejich hodnocení. Právě zde se uplatní znalosti, zkušenosti a všeobecný rozhled zpracovávatele, právě tu se projeví jeho vztah k povaze naměřených dat. Správnou a citlivou volbou regresního modelu lze ze souboru dat vytěžit spoustu informací, naopak zvolením neadekvátního modelu, lze snadno dospět i ke zcela mylným a falešným vývodům. Chcete-li mít v tomto oboru dobré výsledky, pak se musíte obrnit značnou dávkou trpělivosti a již předem počítat s tím, že jen zřídka se vám podaří najít ten správný regresní model hned napoprvé. Z vlastní zkušenosti vím, že k některým modelům se člověk dopravuje až po několika letech marných pokusů.

¹ Váhu jsme povinni akceptovat zejména v případě, že pracujeme se soubory s diametrálně odlišnou kvalitou měření (vizuální a fotoelektrická pozorování jasnosti), nebo tehdy, souvisí-li nejistota jednotlivých měření velikosti závislé veličiny, případně tehdy, nepracujeme-li přímo s naměřenou hodnotou y , ale s její nelineární transformací.

Regresní model představuje množinu podobných funkcí, které se od sebe liší jen různými hodnotami parametrů $\beta_1, \beta_2, \dots, \beta_g$:

$$F(x) = F(\beta_1, \beta_2, \dots, \beta_g, x).$$

Uspořádanou g -ticí parametrů β_j je výhodné zapisovat jako g -rozměrný vektor nebo sloupcovou matici β o rozměrech $g \times 1$ (g řádků a 1 sloupec):

$$\beta' = [\beta_1, \beta_2, \dots, \beta_g]$$

Předpokládejme nyní, že jsme v rámci regresního modelu zvolili nějakou konkrétní hodnotu vektoru parametrů β pro i -té měření $\{x_i, y_i\}$ pak lze vyjádřit odchylku tohoto měření od dané závislosti e_i vztahem:

$$y_i = F(x_i, \beta) + e_i. \quad (2)$$

Je zřejmé, že čím menší budou odchylky, tím lepší bude proložení pozorované závislosti mezi veličinami y a x .

Naším úkolem nyní bude vybrat z třídy funkcí $F(x, \beta)$ popsaných vektorem β , najít takový vektor $\beta = \mathbf{b}$, pro nějž budou odchylky $\{e_i\}$ minimální. Onu podmínku minimálnosti je ovšem třeba nejprve matematicky precizovat.

Nejčastěji používanou, a z mnoha důvodů nejoblíbenější (nikoli však jedinou²), je podmínka, aby součet kvadrátů odchylek pro všechny body měření byl minimální. Z této podmínky pak vychází tzv. metoda nejmenších čtverců které se budeme nadále věnovat.

Zavedme si nejprve skalární veličinu $S(\beta)$, zvanou též součet čtverců odchylek³:

$$S(\beta) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n [y_i - F(x_i, \beta)]^2, \quad S(\beta) = \sum_{i=1}^n e_i^2 w_i = \sum_{i=1}^n [y_i - F(x_i, \beta)]^2 w_i.$$

Nyní hledáme takový vektor β , ($\beta = \mathbf{b}$) pro nějž je součet čtverců odchylek $S(\beta = \mathbf{b})$ minimální.

Ze zadání je zřejmé, že suma čtverců odchylek S musí být nutně veličinou nezápornou. V reálných případech je to navíc veličina kladná, a to ze dvou důvodů: 1) jen zřídka se nám podaří regresní model vybrat natolik dobře, aby pozorovanou závislost popisoval realisticky v celém rozsahu i v detailech, 2) i kdyby se nám to podařilo, pak je nutno počítat s tím, že závislou veličinou y neměříme nikdy absolutně přesně. Každé měření je zatíženo chybou měření, u nichž budeme předpokládat, že odchylky jimi způsobené mají náhodné rozložení.

² Jinou takovou podmínkou může být minimálnost součtu absolutních hodnot odchylek nebo jejich čtvrtých mocnin. Nicméně takto definované podmínky se používají jen zřídka, a ve zcela odůvodněných případech.

³ Obdobně si lze zavést součet čtverců odchylek i v obecnějším případě, kdy závisle proměnná pozorovaná veličina y je funkcí několika nezávislých proměnných (x_1, x_2, \dots, x_m), udávajících vektor \mathbf{x} s m složkami, přičemž každému z měření můžeme přisoudit jistou váhu w . i -té měření je pak dáno uspořádanou $(m+2)$ -ticí čísel $\{x_i, y_i, w_i\}$. Regresní model pak bude funkcí m nezávislých proměnných a g parametrů. Hledáme nyní takovou funkci z regresního modelu, pro niž je funkcionál

$$S(\beta) = \sum_{i=1}^n [y_i - F(\mathbf{x}_i, \beta)]^2 w_i$$

minimální.

Funkci $S(\boldsymbol{\beta})$ si můžete představit jako zprohýbanou plochu v $(g+1)$ rozměrném prostoru, kde g rozměrů je vyhrazeno pro složky vektoru $\boldsymbol{\beta}$ a $(g+1)$ -tý rozměr je rezervován pro funkční hodnotu $S(\boldsymbol{\beta})$. Obecně může mít taková plocha dosti komplikovaný vzhled. Nicméně vždy na ní můžeme jedno nebo i více lokálních minim, z nichž ovšem jen některá budou mít nějaký dobrý fyzikální smysl. Pro vyhledávání minim v průběhu funkce dané několika proměnnými je vypracována řada metod, vesměs numerických. V omezeném počtu případů však lze k výsledku dospět i postupy analytické matematiky.

Fyzikálně reálné minimum se vyznačuje tím, že funkce $S(\boldsymbol{\beta})$ v něm je spojitá a spojitě jsou i všechny parciální derivace, které v bodu minima jsou rovny nule. Platí tedy:

$$\left. \frac{\partial S(\boldsymbol{\beta})}{\partial \beta_k} \right|_{\boldsymbol{\beta}=\mathbf{b}} = \mathbf{0}, \text{ pro všechna } k = (1, 2, \dots, g).$$

Dosadíme-li za $S(\boldsymbol{\beta})$ dostaneme g rovnic ve tvaru:

$$\sum_{i=1}^n \frac{\partial F(x_i, \mathbf{b})}{\partial \beta_k} F(x_i, \mathbf{b}) = \sum_{i=1}^n \frac{\partial F(x_i, \mathbf{b})}{\partial \beta_k} y_i, \quad \sum_{i=1}^n \frac{\partial F(x_i, \mathbf{b})}{\partial \beta_k} F(x_i, \mathbf{b}) w_i = \sum_{i=1}^n \frac{\partial F(x_i, \mathbf{b})}{\partial \beta_k} y_i w_i.$$

Nastává-li pak v bodě $\boldsymbol{\beta} = \mathbf{b}$ minimum, pak je splněno všech g podmínkových rovnic daných výše uvedeným vztahem. Funkce $F(x, \mathbf{b})$, nazývaná též **regresní funkce**, je pak onou hledanou funkcí, která představuje nejlepší přiblížení (nebo je jedním z nich) k průběhu funkční závislosti $y(x)$.

Ještě poznámku. Při hledání extrémů (minima nebo maxima) skalárních funkce je vhodné si zavést pojem gradientu funkce. Gradient v daném bodě je vektor orientovaný v opačném směru než spádnice, přičemž délka vektoru je tím větší, čím strměji v daném bodě funkce probíhá. Číselně jsou složky vektoru gradientu funkce S , která je funkcí g proměnných parametrů, rovny parciálními derivacím podle těchto parametrů:

$$\text{grad } S(\boldsymbol{\beta}) = \left[\frac{\partial S}{\partial \beta_1}, \frac{\partial S}{\partial \beta_2}, \dots, \frac{\partial S}{\partial \beta_g} \right].$$

Gradient lze takto podle potřeby chápat jako buď jako vektor o g složkách nebo řádkovou matici s g sloupci. Pomocí gradientu součtu čtverců odchylek lze podmínku pro nalezení minima funkce lze pak elegantně zapsat:

$$\text{grad } S(\mathbf{b}) = \mathbf{0},$$

kde $\mathbf{0}$ je vektorem o g složkách, jež jsou všechny rovny nule. Podmínka tak říká, že minimum skalární funkce nastává v tom bodě, kdy všechny složky gradientu funkce jsou rovny nule. Velikost vektoru gradientu je v tomto bodě nulová, jsme na dně - hlouběji se již dostat nelze. Popisované metodě hledání minima skalární funkce se proto říká též *gradientní metoda* (gradient method).

Dosadíme-li nyní výraz pro sumu čtverců odchylek dojdeme po jistých úpravách k jediné vektorové podmínce:

$$\sum_{i=1}^n g(x_i, \mathbf{b}) F(x_i, \mathbf{b}) = \sum_{i=1}^n g(x_i, \mathbf{b}) y_i, \quad \sum_{i=1}^n g(x_i, \mathbf{b}) F(x_i, \mathbf{b}) w_i = \sum_{i=1}^n g(x_i, \mathbf{b}) y_i w_i,$$

kde $g(x_i, \mathbf{b})$ je gradientem funkce $F(x_i, \mathbf{b})$ v bodu $x_i, \boldsymbol{\beta} = \mathbf{b}$.

Soustavu g rovnic o g neznámých (b_j) pak lze standardním způsobem řešit. Nalezením všech hledaných koeficientů je pak nalezena i regresní funkce, kde $\boldsymbol{\beta} = \mathbf{b}$. Pokud nás nezajímá přesnost měření, hodnověrnost proložení, chyby parametrů a neurčitost předpovědi, pak jsme hotovi. V opačném případě budeme postupovat dále.

Hodnoty V_{kj} definují prvky čtvercové symetrické matice \mathbf{V} rozměru $g \times g$, zatímco hodnoty pravých stran rovnice U_k definují prvky vektoru \mathbf{U} s g prvky (sloupcová matice $g \times 1$). Nyní lze celou soustavu rovnic zapsat ještě elegantněji:

$$\mathbf{V} \mathbf{b} = \mathbf{U}.$$

Definujme si nyní tzv. *kovarianční matici* \mathbf{H} . Jde o čtvercovou matici rozměru $g \times g$, která je matricí inverzní k matici \mathbf{V} . Platí tedy o ní:

$$\mathbf{H} = \mathbf{V}^{-1}, \quad \mathbf{H} \mathbf{V} = \mathbf{V} \mathbf{H} = \mathbf{I},$$

kde \mathbf{I} je jednotková matice. Vynásobím-li zleva obě dvě strany rovnice maticí \mathbf{H} , dostanu přímý vztah pro hledanou vektorovou matici \mathbf{b} :

$$\mathbf{b} = \mathbf{H} \mathbf{U}.$$

Pro další výpočty je výhodné pracovat s vektorovou funkcí $\mathbf{g}(x)$, která je gradientem funkce regresního modelu $F(x, \boldsymbol{\beta})$. V případě, že je regresní model lineární kombinací funkcí $f_1(x), f_2(x), \dots, f_g(x)$, je vyjádření gradientu velmi prosté:

$$\mathbf{g}(x) = \text{grad } F(x, \boldsymbol{\beta}) = [f_1(x), f_2(x), \dots, f_g(x)].$$

Funkční hodnota regresní funkce $F(x, \mathbf{b})$ je současně i předpovědí $y_p(x)$ pro zvolenou hodnotu x .

$$y_p(x) = F(x, \boldsymbol{\beta}) = \sum_{j=1}^g b_j f_j(x) = \mathbf{g}(x) \mathbf{b}.$$

K odhadu nejistoty proložení a chyb určení jednotlivých parametrů a předpovědi je nutno nejdříve vypočítat hodnotu zbytkového (reziduálního) součtu čtverců odchylek R pro nalezenou hodnotu $\boldsymbol{\beta} = \mathbf{b}$, kdy je tento součet minimální:

$$R = S(\mathbf{b}) = \sum_{i=1}^n [y_i - y_p(x_i)]^2 = \sum_{i=1}^n [y_i - \mathbf{g}(x_i) \mathbf{b}]^2,$$

$$R = S(\mathbf{b}) = \sum_{i=1}^n w_i [y_i - y_p(x_i)]^2 = \sum_{i=1}^n w_i [y_i - \mathbf{g}(x_i) \mathbf{b}]^2.$$

Pomocí reziduální součtu čtverců odchylek R lze odhadnout velikost *střední kvadratické odchylky* jednoho měření σ . Ta s počtem měření n , počtem stupňů volnosti g a součtem čtverců odchylek R souvisí takto:

$$\sigma = \sqrt{\frac{R}{n-g}}, \quad \sigma = \sqrt{\frac{R}{\bar{w}(n-g)}}.$$

Odhad *nejistoty* (uncertainty) určení velikosti k -tého parametru vektoru \mathbf{b} , $\delta \mathbf{b}$, je dán vztahem:

$$\delta \mathbf{b} = \sigma \sqrt{\text{diag}(\mathbf{H})} \quad \delta \mathbf{b} = \sigma \sqrt{\bar{w} \text{diag}(\mathbf{H})}.$$

Odhad nejistoty funkční hodnoty nalezené regresní funkce $y_p(x)$ v daném bodě x , neboli *předpovědi* v bodě x , $\delta y_p(x)$, je dána vztahem:

$$\delta y_p(x) = \sigma \sqrt{\mathbf{g}(x) \mathbf{H} \mathbf{g}'(x)} \quad \delta y_p(x) = \sigma \sqrt{\bar{w} \mathbf{g}(x) \mathbf{H} \mathbf{g}'(x)}.$$

Po určitých úpravách vztahu pro součet čtverců odchylek pro případ lineární regrese, lze veličinu $S(\boldsymbol{\beta})$ uvést v instruktivním tvaru:

$$S(\boldsymbol{\beta}) = R + \sum_{j=1}^g (\beta_j - b_j)^2 \sum_{i=1}^n f_j^2(x_i), \quad S(\boldsymbol{\beta}) = R + \sum_{j=1}^g (\beta_j - b_j)^2 \sum_{i=1}^n w_i f_j^2(x_i).$$

Ze zápisu je okamžitě patrné, že funkce S má tvar paraboloidu s minimem o hodnotě v bodu $\boldsymbol{\beta} = \mathbf{b}$. Má tedy jediné a tudíž absolutní minimum.

Velmi elegantní lze lineární regresi řešit použitím maticového počtu. Definujme si tři matice \mathbf{X} , \mathbf{Y} , případně \mathbf{W} o rozměrech postupně $n \times g$, $n \times 1$ a diagonální $n \times n$:

$$\mathbf{X} = \begin{bmatrix} \mathbf{g}(x_1) \\ \mathbf{g}(x_2) \\ \vdots \\ \mathbf{g}(x_n) \end{bmatrix} = \begin{bmatrix} f_1(x_1) & f_2(x_1) & \dots & f_g(x_1) \\ f_1(x_2) & f_2(x_2) & \dots & f_g(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ f_1(x_n) & f_2(x_n) & \dots & f_g(x_n) \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{W} = \text{diag} \left(\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \right).$$

$$\mathbf{V} = \mathbf{X}' \mathbf{X}, \quad \mathbf{V} = \mathbf{X}' \mathbf{W} \mathbf{X},$$

$$\mathbf{H} = \mathbf{V}^{-1} = (\mathbf{X}' \mathbf{X})^{-1}, \quad \mathbf{H} = \mathbf{V}^{-1} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1}.$$

$$\mathbf{U} = \mathbf{X}' \mathbf{Y}, \quad \mathbf{U} = \mathbf{X}' \mathbf{W} \mathbf{Y}.$$

$$\mathbf{b} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} = (\mathbf{U}' / \mathbf{V})' = \mathbf{X} \setminus \mathbf{Y}, \quad \mathbf{b} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{Y} = (\mathbf{U}' / \mathbf{V})'.$$

Definujme nyní vektorovou sloupcovou matici předpovědi $\mathbf{Y}_p = [y_p(x_1); y_p(x_2); \dots; y_p(x_n)]$

$$\mathbf{Y}_p = \mathbf{X} \mathbf{b}.$$

Reziduální součet čtverců odchylek je pak dán vztahem:

$$R = (\mathbf{Y} - \mathbf{Y}_p)' (\mathbf{Y} - \mathbf{Y}_p) = \mathbf{Y}' \mathbf{Y} - \mathbf{b}' \mathbf{U}, \quad R = (\mathbf{Y} - \mathbf{Y}_p)' \mathbf{W} (\mathbf{Y} - \mathbf{Y}_p) = \mathbf{Y}' \mathbf{W} \mathbf{Y} - \mathbf{b}' \mathbf{U}.$$

$$\delta \mathbf{Y}_p = \sqrt{\frac{R}{n-g} \text{diag}(\mathbf{X} \mathbf{H} \mathbf{X}')}$$

3 Základní regresní modely - Aplikace MNČ

Následuje několik praktických příkladů aplikace prosté i obecné metody nejmenších čtverců, které mají ilustrovat způsob, jak se má MNČ používat. Pokud tyto příklady někomu připadnou jako triviální, pak se nemýlí, neboť jde o záměr.

V řadě příkladů budou s výhodou použity některé střední veličiny, které zde zavedeme. Aritmetický průměr veličin x a y :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \bar{x} = \frac{1}{S_w} \sum_{i=1}^n x_i w_i, \quad \bar{y} = \frac{1}{S_w} \sum_{i=1}^n y_i w_i,$$

Dále střední hodnoty součinu veličin x a y v různých mocninách:

$$\overline{x^m y^l} = \frac{1}{n} \sum_{i=1}^n x_i^m y_i^l, \quad \overline{x^m y^l} = \frac{1}{S_w} \sum_{i=1}^n x_i^m y_i^l w_i,$$

kde m a l jsou celá nezáporná čísla $0, 1, \dots$

Užitečné je též zavedení tzv. rozptylu u_{xx} a u_{yy} a směrodatných odchylek s_x a s_y souboru veličin x a y a míru korelace mezi nimi u_{xy} .

$$u_{xx} = \overline{x^2} - \bar{x}^2, \quad s_x = \sqrt{u_{xx}}, \quad u_{yy} = \overline{y^2} - \bar{y}^2, \quad s_y = \sqrt{u_{yy}}, \quad u_{xy} = \overline{xy} - \bar{x}\bar{y},$$

Bezrozměrný koeficient korelace r :

$$r = \frac{\overline{xy} - \bar{x}\bar{y}}{s_x s_y} = \sqrt{\frac{u_{xy}^2}{u_{xx} u_{yy}}} = \frac{u_{xy}}{s_x s_y}.$$

Lze ukázat, že koeficient korelace r nabývá hodnot mezi -1 a 1 , přičemž 0 je roven tehdy, kdy mezi veličinami x a y neexistuje žádná lineární korelace, ± 1 je roven tehdy, kdy jsou všechny hodnoty $\{x_i, y_i\}$ uloženy na jediné přímce.

3.1 Střední hodnota veličin se stejnou vahou (PMNČ)

Nasvědčuje-li měření n dvojic $\{x_i, y_i\}$ tomu, že mezi x a y neexistuje žádná závislost a že hodnota $y(x)$ je v mezích chyb nejspíš konstantní, postavíme regresní model takto:

$$y_i = \beta + e_i.$$

Optimální hodnotu β , při níž je suma kvadrátů odchylek e_i minimální, b , nazveme střední hodnotou. Najdeme ji minimalizací funkcionálu $S(\beta)$:

$$S(\beta) = \sum_{i=1}^n (y_i - \beta)^2 = \sum_{i=1}^n y_i^2 - 2\beta \sum_{i=1}^n y_i + n\beta^2 = n(\overline{y^2} - 2\beta \bar{y} + \beta^2).$$

Grafem funkce je parabola s minimem v bodu $\beta = \bar{y}$, přičemž s minimem $R = S(\beta = \bar{y})$:

$$R = n(\overline{y^2} - \bar{y}^2) = n u_{yy}.$$

I když minimalizaci funkce $S(\beta)$ lze vypočítat přímo, zkusme si nyní ze cvičných důvodů všechny potřebné vztahy odvodit pomocí maticových vztahů.

$$\mathbf{X} = [1; \dots; 1], \quad \mathbf{Y} = [y_1; y_2; y_3; \dots; y_n],$$

$$\mathbf{V} = \mathbf{X}'\mathbf{X} = n, \quad \mathbf{H} = (\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{n}, \quad \mathbf{U} = \mathbf{X}'\mathbf{Y} = \sum_{i=1}^n y_i = n \bar{y}.$$

Jak patrně, nikde se v důležitých veličinách nevyskytují veličiny x_i , tudíž na nich nezáleží a mohou nabývat libovolnou hodnotu.

$$\mathbf{g}(x) = 1, \quad \mathbf{b} = \mathbf{H}\mathbf{U} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = (\mathbf{U}'/\mathbf{V})' = \frac{1}{n}n\bar{y} = \bar{y}.$$

Střední hodnota podle MNČ je tedy přímo rovna aritmetickému průměru.

$$y_p(x) = \mathbf{g}(x)\mathbf{b} = \bar{y}, \quad R = \mathbf{Y}'\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{Y} = \sum_{i=1}^n y_i^2 - \bar{y}n\bar{y} = n u_{yy}.$$

$$\sigma = \sqrt{\frac{R}{n-1}} = s_y \sqrt{\frac{n}{n-1}}, \quad \delta b = \sigma \sqrt{\text{diag}(\mathbf{H})} = \frac{\sigma}{\sqrt{n}}, \quad \delta y_p = \frac{\sigma}{\sqrt{n}}.$$

3.2 Střední hodnota veličin s nesterjnou váhou (OMNČ)

Stanovení střední hodnoty veličin s nesterjnou váhou je nejjednodušší úlohou řešitelnou OMNČ. Předpokládejme, že máme n trojic veličin $\{x_i, y_i, w_i\}$, kde w_i je váha i -tého měření veličiny y_i , a nezávislá veličina x_i , na jejíž velikosti však v tomto případě nijak nezáleží. Regresní model bude tž jako v případě A.

$$y_i = \beta + e_i.$$

Optimální hodnotu β , při níž je suma váhovaných kvadrátů odchylek e_i minimální, b , nazveme střední hodnotou. Najdeme ji minimalizací funkce $S(\beta)$:

$$S(\beta) = \sum_{i=1}^n w_i (y_i - \beta)^2.$$

I když minimalizace je i v tomto případě jednoduchá operace, opět využijeme maticových vztahů.

$$\mathbf{X} = \text{ones}(n,1), \quad \mathbf{Y} = [y_1; y_2; y_3; \dots; y_n], \quad \mathbf{W} = \text{diag}[w_1, w_2, \dots, w_n].$$

$$\mathbf{X}'\mathbf{W} = [w_1, w_2, \dots, w_n], \quad \mathbf{Y}'\mathbf{W} = [y_1 w_1, y_2 w_2, \dots, y_n w_n]$$

$$\mathbf{V} = \mathbf{X}'\mathbf{W}\mathbf{X} = \sum_{i=1}^n w_i = n\bar{w}, \quad \mathbf{H} = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1} = \frac{1}{\bar{w}}, \quad \mathbf{U} = \mathbf{X}'\mathbf{W}\mathbf{Y} = n\bar{y}$$

$$\mathbf{g}(x) = 1, \quad \mathbf{b} = \mathbf{H}\mathbf{U} = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{Y} = \bar{y}.$$

Střední hodnota podle OMNČ je tedy přímo rovna váhovanému aritmetickému průměru.

$$y_p(x) = \mathbf{F}(x)\mathbf{b} = \bar{y}, \quad R = \mathbf{Y}'\mathbf{W}\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{W}\mathbf{Y} = \sum_{i=1}^n y_i^2 w_i - \bar{y} \sum_{i=1}^n y_i w_i = n \bar{w} u_{yy}$$

$$\sigma = \sqrt{\frac{R}{(n-1)\bar{w}}} = s_y \sqrt{\frac{n}{n-1}}, \quad \delta b = \sigma \sqrt{\bar{w}\mathbf{H}} = \frac{\sigma}{\sqrt{n}}, \quad \delta y_p = \sigma \sqrt{\bar{w}[\mathbf{g}(x)\mathbf{H}\mathbf{g}'(x)]} = \frac{\sigma}{\sqrt{n}}.$$

Za povšimnutí jistě stojí, že vztahy pro b , s , δb a δy_p jsou formálně stejné jako v případě A, kdy se tentýž problém řešil pro stejné váhy. Rozdíl ovšem je v tom, jak jsou definovány střední veličiny, z nichž se při výpočtu vychází.

3.3 Přímka jdoucí počátkem I (PMNČ)

Občas se můžeme setkat se situací, kdy je jeden nebo více bodů závislosti pevně fixováno. Z této skutečnosti musíme při volbě regresního modelu vycházet. Nejjednodušším příkladem toho druhu je naše očekávání, že n bodů o souřadnicích $[x_i, y_i]$ se stejnými váhami lze proložit přímkou jdoucí bodem o souřadnicích $[0, 0]$, neboli počátkem. Regresní model je pak:

$$y_i = \beta x_i + e_i$$

Optimální hodnotu β , při níž je suma kvadrátů odchylek e_i minimální, b , nazveme tentokrát středním koeficientem úměrnosti.

$$\mathbf{X} = [x_1; x_2; x_3; \dots; x_n] \quad \mathbf{Y} = [y_1; y_2; y_3; \dots; y_n]$$

$$V = \mathbf{X}'\mathbf{X} = \sum_{i=1}^n x_i^2 = n\overline{x^2}, \quad H = (\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{n\overline{x^2}}, \quad U = \mathbf{X}'\mathbf{Y} = \sum_{i=1}^n x_i y_i = n\overline{xy}, \quad \mathbf{g}(x) = x,$$

$$b = \left(\frac{\mathbf{U}'}{\mathbf{V}} \right)' = \frac{U}{V} = \frac{\overline{xy}}{\overline{x^2}}, \quad y_p(x) = b x, \quad R = \mathbf{Y}'\mathbf{Y} - b \mathbf{X}'\mathbf{Y} = n(\overline{y^2} - b\overline{xy}) = n \left[\overline{y^2} - \frac{(\overline{xy})^2}{\overline{x^2}} \right],$$

$$\sigma = \sqrt{\frac{R}{n-1}} = \sqrt{\frac{n}{n-1}(\overline{y^2} - b\overline{xy})} = \sqrt{\frac{n}{n-1}(\overline{y^2} - b^2\overline{x^2})} = \sqrt{\frac{n}{n-1}(\overline{y^2} - b\overline{xy})} = \sqrt{\frac{n}{n-1} \left[\overline{y^2} - \frac{(\overline{xy})^2}{\overline{x^2}} \right]},$$

$$\delta b = \sigma \sqrt{H} = \frac{\sigma}{\sqrt{n\overline{x^2}}} = \sqrt{\frac{1}{n-1} \left(\frac{\overline{y^2}}{\overline{x^2}} - b^2 \right)}, \quad \delta y_p = \sigma \sqrt{\mathbf{g}(x) \mathbf{H} \mathbf{g}'(x)} = \sigma \sqrt{\frac{x^2}{n\overline{x^2}}}.$$

Poznámka: Pokusme vypočítat koeficient úměrnosti jinak. Uvažme, že pomocí každé z dvojic x_i, y_i lze vypočítat „individuální“ koeficient úměrnosti b_i : $b_i = y_i/x_i$ a střední koeficient úměrnosti, který si nyní označíme b' , by pak logicky měl být roven aritmetickému průměru jednotlivých b_i :

$$b' = \frac{1}{n} \sum b_i = \frac{1}{n} \sum \frac{y_i}{x_i}$$

Jak patrně, tento vztah se od výše uvedeného vztahu pro hodnotu středního koeficientu úměrnosti liší a žádnou z dovolených matematických operací nelze tyto dva vztahy ztotožnit. Jak to vysvětlit?

Vysvětlení plyne z předpokladu, na němž je prostá MNČ postavena: rozptyl měření y od reálného průběhu daného funkční závislosti $y(x)$ má povahu náhodné veličiny a zejména nijak nezávisí na hodnotě další měřené veličiny x . Je-li tato podmínka splněna pro množinu měření $\{x_i, y_i\}$, pak ovšem nemůže být splněna pro veličinu y_i/x_i , jejíž očekávaná nepřesnost je nepřímo úměrná hodnotě x_i . Různě velkou očekávanou nepřesnost je třeba v OMNČ ocenit různým ováhováním bodů, kdy váha jednotlivého bodu - w_i - bude nepřímo úměrná kvadrátu očekávaného rozptylu měření, čili v tomto případě by měla být úměrná x_i^2 . Položme proto přímo, že $w_i = x_i^2$. Aritmetický průměr z množiny b_i se započítáním jejich individuálních vah se vypočte podle vztahu:

$$\bar{b} = \frac{\sum_{i=1}^n b_i w_i}{\sum_{i=1}^n w_i} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{\overline{xy}}{\overline{x^2}} = b.$$

3.4 Příмка jdoucí počátkem II (OMNČ)

Oproti případu C budeme navíc předpokládat, že každému z bodů měření o souřadnicích $[x_i, y_i]$ bude přisouzena určitá individuální váha w_i . Přesvědčete se sami, že pokud budeme počítat se středními váženými veličinami a jejich součiny, pak vztahy pro střední koeficient úměrnosti b , střední váženou kvadratickou odchylku jednoho měření σ , chyba koeficientu δb a chyba předpovědi δy_p , budou formálně stejné jako v případu 3.3.

3.5 Obecná příмка (OMNČ)

Snad nejběžnější úlohou, s níž se při zpracování pozorování můžeme setkat je, jak zjistit parametry předpokládané lineární závislosti mezi veličinami y a x , nebo jak proložit body v grafu přímkou. Regresní model je zřejmý:

$$y_i = \alpha + \beta x_i + e_i.$$

Příмка necht' je prokládána n body o souřadnicích $[x_i, y_i]$, přičemž každému z bodů je přisouzena jeho individuální váha w_i . Řešením úlohy je nalezení takové dvojice parametrů a a b , pro něž je suma váhovaných čtverců odchylek $S(\alpha, \beta)$ minimální:

$$S(\alpha, \beta) = \sum_{i=1}^n w_i (y_i - \alpha + \beta x_i)^2.$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & w_n \end{bmatrix}.$$

$$\mathbf{V} = \mathbf{X}' \mathbf{W} \mathbf{X} = n \bar{w} \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \overline{x^2} \end{bmatrix}, \quad \mathbf{U} = \mathbf{X}' \mathbf{W} \mathbf{Y} = n \bar{w} \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix}, \quad \mathbf{H} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} = \frac{1}{n \bar{w} u_{xx}} \begin{bmatrix} \overline{x^2} & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix},$$

$$\mathbf{b} = \mathbf{H} \mathbf{U} = \begin{bmatrix} a \\ b \end{bmatrix} = \frac{1}{u_{xx}} \begin{bmatrix} \overline{x^2} & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix} = \frac{1}{u_{xx}} \begin{bmatrix} \overline{x^2} \bar{y} - \bar{x} \overline{xy} \\ -\bar{x} \bar{y} + \overline{xy} \end{bmatrix}, \Rightarrow a = \frac{\overline{x^2} \bar{y} - \bar{x} \overline{xy}}{u_{xx}}, \quad b = \frac{u_{xy}}{u_{xx}} = r \frac{s_y}{s_x},$$

kde r je korelační koeficient.

$$\mathbf{g}(x) = [1 \ x], \quad y_p(x) = \mathbf{g}(x) \mathbf{b} = [1 \ x] \begin{bmatrix} a \\ b \end{bmatrix} = a + b x.$$

Přesvědčte se, že platí: $y_p(\bar{x}) = \bar{y}$, což jinými slovy znamená, že regresní přímka prochází těžištěm.

$$R = \mathbf{Y}' \mathbf{W} \mathbf{Y} - \mathbf{b}' \mathbf{X}' \mathbf{W} \mathbf{Y} = n \bar{w} (\bar{y}^2 - a \bar{y} - b \bar{xy}), \quad \sigma = \sqrt{\frac{R}{(n-g)\bar{w}}} = \sqrt{\frac{n(\bar{y}^2 - a \bar{y} - b \bar{xy})}{n-2}},$$

$$\delta b = \sigma \sqrt{\bar{w} \mathbf{H}_{22}} = \frac{1}{\sqrt{n}} \frac{\sigma}{s_x}, \quad \delta a = \sigma \sqrt{\bar{w} \mathbf{H}_{11}} = \frac{\sigma}{\sqrt{n}} \frac{\sqrt{\bar{x}^2}}{s_x} = \delta b \sqrt{\bar{x}^2}.$$

Nejistota směrnice přímky b tedy nezávisí na umístění počátku, zatímco chyba absolutního členu a ano. Minimální je tato chyba v případě, kdy počátek souřadnic ztotožníme s těžištěm. Chyba pak bude $\delta a = \sigma/\sqrt{n}$.

Chyba předpovědi funkční hodnoty y_p v bodě x - δy_p - je dána vztahem:

$$\delta y_p(x) = \sigma \sqrt{\bar{w} [\mathbf{g}(x) \mathbf{H} \mathbf{g}'(x)]} = \frac{\sigma}{\sqrt{n}} \sqrt{1 + \frac{(x - \bar{x})^2}{s_x^2}}.$$

Vidíme, že podle očekávání je chyba předpovědi minimální v oblasti v těsné blízkosti těžiště, ve velkých vzdálenostech od něj je asymptoticky přímo úměrná této vzdálenosti.

Absolutní člen a lze geometricky interpretovat jako úsek na ose y , který na ní vytíná regresní přímka. Neurčitost polohy tohoto průsečíku udává chyba předpovědi δy_p v bodě $x = 0$. Číselně tato chyba je rovna chybě absolutního členu, tak jak jej udává výše uvedený vztah pro δa .

Poznámka:

Doporučuji ještě před výpočtem provést transformaci závislé proměnné tak, že je budeme vztahovat k jejich střední vážené hodnotě: $t_i = x_i - \bar{x}$. Transformace ovlivní jen absolutní člen $a' = \bar{y}$, $\delta a' = \sigma/\sqrt{n}$,

$$y_p(x) = \bar{y} + \bar{b}t, \quad \delta y_p(x) = \frac{\sigma}{\sqrt{n}} \sqrt{1 + \frac{t^2}{s_x^2}}.$$

Tímto postupem se výrazně omezí vliv zaokrouhlovacích chyb při výpočtu a zvýší se tak přesnost a spolehlivost výsledku.

4 Úloha

Proložte přímku těmito daty $\{x, y, w\}$ a vypočtěte její parametry:

| | | | | | | | | | | | |
|------|-------|---|------|------|---|------|------|---|------|------|---|
| 0.11 | -0.06 | 1 | 0.82 | 0.82 | 2 | 0.36 | 0.38 | 3 | 0.57 | 0.18 | 1 |
| 0.66 | 0.69 | 2 | 0.67 | 0.78 | 2 | 0.55 | 0.92 | 2 | 0.70 | 0.36 | 2 |
| 0.37 | 0.34 | 3 | 1.00 | 0.86 | 1 | 0.26 | 0.21 | 1 | 0.96 | 0.85 | 2 |
| 0.14 | -0.09 | 4 | 0.96 | 0.83 | 3 | 0.60 | 1.04 | 1 | 0.75 | 0.71 | 3 |
| 0.57 | 0.80 | 5 | 0.06 | 0.12 | 4 | 0.05 | 0.35 | 4 | 0.74 | 0.74 | 4 |

- Tak aby procházela počátkem
- Obecnou přímku, diskutujte přitom, zda není model a) lepší

- c) Zaměňte závislou a nezávislou proměnnou a výsledek porovnejte. Dokažte, že poměr směrnic za těchto okolností je roven r^2 !
- d) Řešte vše pro situaci s vahami a bez nich. Výsledky porovnejte.
- e) Zkuste proložit závislost polynomem vyššího stupně a diskutujte oprávněnost toho regresního modelu.
- f) Vykreslete graf s body, proloženou přímkou a nejistotou proložení ($y_p \pm \delta y_p$)

5 Řešení úlohy

a) Přímka, jež prochází počátkem:

$$s = 0.216; y = (0,97 \pm 0,08) x$$

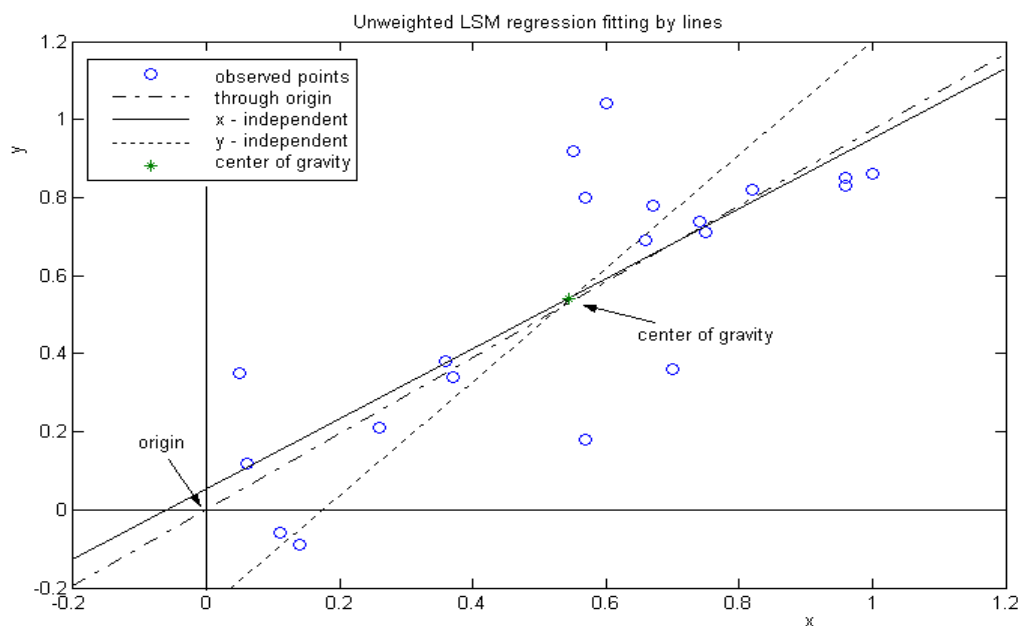
b) Obecná přímka, x nezávislá proměnná:

$$s = 0.220; y = (0,05 \pm 0,11) + (0,90 \pm 0,17) x, r = 0,78. \text{ , model a) je zřejmě lepší.}$$

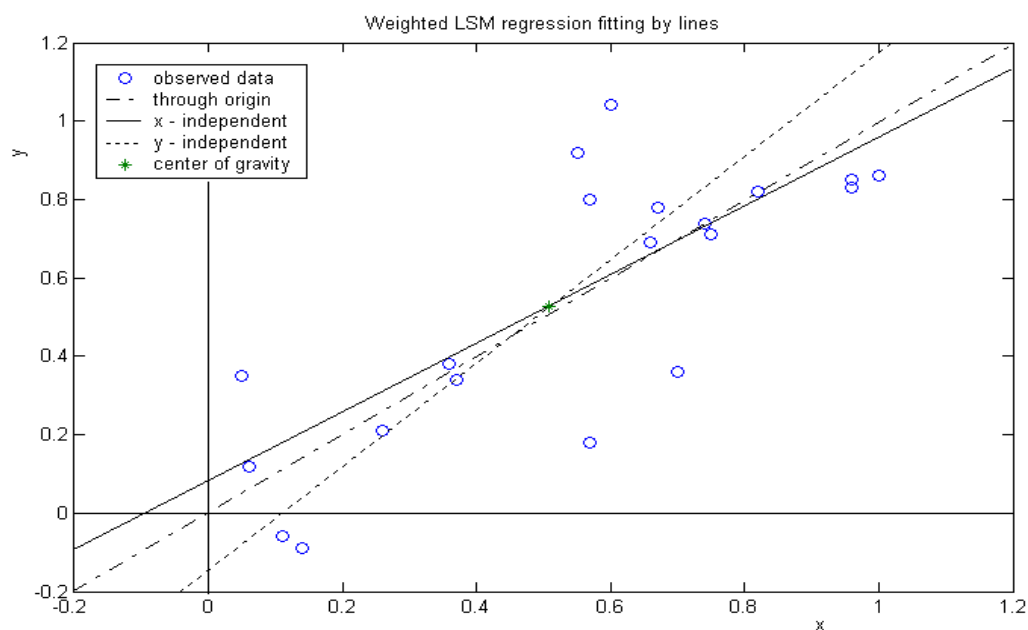
c) Obecná přímka, x nezávislá proměnná:

$$s = 0.193; x = (0,17 \pm 0,08) + (0,69 \pm 0,13) x, r = 0,78.$$

Odmocnina poměru směrnic $\sqrt{0,90 \cdot 0,69} = 0,78$ je rovna r .



d) Situace s vahami



Přímka, jež prochází počátkem:

$$s = 0.198; y = (0,99 \pm 0,08) x$$

Obecná přímka, x nezávislá proměnná:

$$s = 0.198; y = (0,05 \pm 0,11) + (0,90 \pm 0,17) x, r = 0,82. , \text{ model a) je stejně dobrý.}$$

Obecná přímka, x nezávislá proměnná:

$$s = 0.185; x = (0,11 \pm 0,08) + (0,76 \pm 0,13) x, r = 0,82.$$

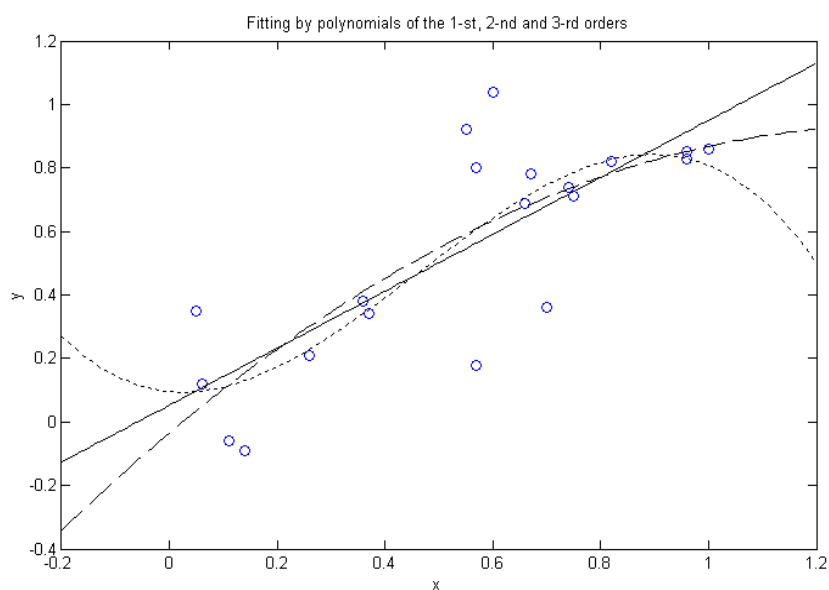
Odmocnina poměru směrníc $\sqrt{0,90 \cdot 0,76} = 0,82$ je rovna r .

e) Proložení polynomem 2. stupně (x je nezávisle proměnná, neváženo): $s = 0,222$,

$$y = (-0,04 \pm 0,14) + (1,43 \pm 0,63) x - (0,53 \pm 0,60) x^2;$$

proložení polynomem 3. stupně: $s = 0,224$

$$y = (0,10 \pm 0,21) + (-0,15 \pm 1,97) x + (3,1 \pm 4,4) x^2 - (2,3 \pm 2,7) x^3$$



f) Vykreslete graf s body, proloženou přímkou a nejistotou proložení ($y_p \pm \delta y_p$)

