

Linearna regresija

Predstavljajmo si, izvajamo fizikalni poskus. Imamo vzmet, na kateri merimo raztezek v odvisnosti od sile, s katero napenjamo vzmet (na vzmet obešamo uteži z znano težo in merimo raztezek).

F [N]	1	2	3	4	5
x [cm]	10	20	35	55	80

Velja Hookov zakon, $F = kx$, torej sta x in F med seboj linearno odvisna. Poskušali bomo najti premico, ki se bo najboljše prilegala danim točkam v ravnini.

Če imamo 2 točki $T_1(x_1, y_1), T_2(x_2, y_2)$, lahko brez težav najdemo premico, ki gre točno skozi niju. Če pa je točk več, ni nujno, da obstaja premica, ki gre skozi vse točke. Iščemo taka a in b , da se bo premica kar najboljše prilegala danim točkam. Če gre premica skozi neko točko $T_i(x_i, y_i)$, potem velja:

$$0 = ax_i + b - y_i$$

Ker pa naša premica ne poteka direktno skozi vse točke pride do napake, ki jo bomo v točki T_i označili z ε_i .

$$\varepsilon_i = ax_i + b - y_i$$

Napaka je lahko pozitivna ali negativna, odvisno ali točka leži pod ali nad premico. Želimo zmanjšati velikost vseh napak, ne glede na to, ali so pozitivne ali negativne. (Lahko bi preprosto sešteli absolutne vrednosti napak, ampak pozneje funkcije ne bi mogli odvajati.) Zato seštejemo kvadrate vseh napak.

$$\varepsilon = \sum_{i=1}^N \varepsilon_i^2$$

¹Podatki pridobljeni iz: <http://www.clemson.edu/ces/phoenix/labs/124/shm/>

$$\varepsilon = \sum_{i=1}^N (ax_i - b - y_i)^2$$

Naš cilj je, da minimiziramo to napako. Minimum funkcije pa najdemo tako, da izračunamo, kdaj je odvod funkcije enak 0.

Odводи

Kaj so odvodi? Geometrijski pomen, ničle odvoda?

Delni odvodi

Funkcija je odvisna od večih spremenljivk, zato jo moramo odvajati za vsako spremenljivko posebej. Ker nam odvod ene spremenljivke pove kako se funkcija obnaša samo po tej spremenljivki, temu rečemo delni odvod.

$$\frac{\partial \varepsilon}{\partial a} = \sum_{i=1}^N 2 \frac{\partial (ax_i - b - y_i)}{\partial a} = 2 \sum_{i=1}^N (ax_i - b - y_i) x_i$$

$$\frac{\partial \varepsilon}{\partial b} = \sum_{i=1}^N 2 \frac{\partial (ax_i - b - y_i)}{\partial b} = 2 \sum_{i=1}^N (ax_i - b - y_i) (-1) = -2 \sum_{i=1}^N (ax_i - b - y_i)$$

Rešitev najdemo tako, da ugotovimo, pri katerih a in b sta odvoda enaka 0. To lahko rešimo z uporabo matrik, vendar tega ne zamo, zato se bomo tega lotili s primitvno metodo gradientnega spusta.

Gradientni spust

Gradient nam pove smer največjega naraščanja funkcije. Gradient funkcije $f(\mathbf{x})$ je vektor definiran kot:

$$\nabla f = \left(\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \frac{\partial f}{\partial x_3} \quad \dots \quad \frac{\partial f}{\partial x_n} \right)$$

V našem primeru linearne regresije je gradientni spust:

$$\nabla \varepsilon = \begin{pmatrix} \frac{\partial \varepsilon}{\partial a} & \frac{\partial \varepsilon}{\partial b} \end{pmatrix}$$

Ker iščemo minimum funkcije se bomo premikali v nasprotno smer gradienta. Spremembo spremenljivke a lahko zapišemo kot:

$$\Delta a = -(\nabla \varepsilon)_1 \cdot \lambda$$

kjer λ predstavlja velikost premika.

Naše nove spremenljivke so:

$$a = a + \Delta a$$

$$b = b + \Delta b$$

S tem smo se premaknili za en korak bližje minimumu funkcije, to pomeni parametram a in b pri katerih se bo premica najbolj prilegala našim točkam. Ta postopek ponovimo še n -krat. Večkrat ko ga ponovimo, bolj natančen bo naš rezultat.

Natančnost našega rezultata je odvisna tudi od λ (velikost premika) in začetnih vrednost spremenljivk a in b . Z manjšo velikostjo premika bo rezultat bolj natančen, vendar bomo morali postopek večkrat ponoviti. Ker s postopkom iščemo samo lokalne minimume, nam začetna vrednost spremenljivk določi kateri minimum bomo našli.

Linerana regresija elipse

Imamo N točk krožnice nekega planeta, za katere želimo najti krožnico (elipso), ki se tem točkam najbolj prilega.

$$T_1(x_1, y_1), T_2(x_2, y_2) \dots, T_n(x_n, y_n)$$

Zapišemo splošno enačbo za krivulje 2. reda:

$$Ax^2 + Bxy + Cy^2 + Dx + Ey + F = 0$$

Ker je točk preveč, ne moremo najti elipse, na kateri bodo ležale vse točke. Zato poskušamo najti stožnico (elipso), ki se točkam najboljše prilaga. Za točko T_i velja:

$$\varepsilon_i = Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F$$

Pri čemer je ε_i napaka v tej točki. Iščemo najboljše parametre A, B, C, D, E in F , tako da bo skupna napaka čim manjša.

Ko izračunamo napako za vsako točko dobimo vektor napak. Zanima nas skupna velikost napake, kar je 'dolžina' tega vektorja.

$$\varepsilon = \sqrt{\sum_{i=1}^N (Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)^2}$$

Želimo najti minimum te funkcije, ki je odvisna od 6-ih spremenljivk. Če iščemo minimum te funkcije, je enako kot da bi iskali minimum ε^2 .

$$\varepsilon^2(A, B, C, D, E, F) = \sum_{i=1}^N (Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)^2$$

Funkcija je odvisna od večih spremenljivk, zato jo moramo odvajati za vsako spremenljivko posebej. Ker nam odvod ene spremenljivke pove kako se funkcija obnaša samo po tej spremenljivki, temu rečemo delni odvod.

Izračunamo delne odvode za to funkcijo:

$$\begin{aligned}\frac{\partial \varepsilon}{\partial A} &= \sum_{i=1}^N 2(Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)(x_i^2) \\ \frac{\partial \varepsilon}{\partial B} &= \sum_{i=1}^N 2(Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)(x_iy_i) \\ \frac{\partial \varepsilon}{\partial C} &= \sum_{i=1}^N 2(Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)(y_i^2) \\ \frac{\partial \varepsilon}{\partial D} &= \sum_{i=1}^N 2(Ax_i^2 + Bx_iy_i + Cy_i^2 + Dx_i + Ey_i + F)(x_i)\end{aligned}$$

$$\frac{\partial \varepsilon}{\partial E} = \sum_{i=1}^N 2(Ax_i^2 + Bx_i y_i + Cy_i^2 + Dx_i + Ey_i + F)(y_i)$$

$$\frac{\partial \varepsilon}{\partial F} = \sum_{i=1}^N 2(Ax_i^2 + Bx_i y_i + Cy_i^2 + Dx_i + Ey_i + F)$$

Da najdemo rešitev, moramo ugotoviti kdaj so vsi odvodi enaki 0. Tega ne znamo (možno je z uporabo matrik, ampak ne znamo), zato se bomo reševanja lotili z gradientnim spustom.

Logistična regresija

Glavna razlika med logistično in linearno regresijo je, da pri linearni regresiji določamo zvezno spremenljivko (y je odvisen od x), medtem ko nam pri logistični regresiji model vrne kakšna je verjetnost, da vhodni podatki sodijo v določeno kategorijo. Zato se logistična regresija uporablja za klasifikacijo (npr. prepoznavanje števk idr.).

Dodatno

Ekliptični koordinatni sistem

Eden izmed koordinatnih sistemov za določanje lege nebesnih teles. Lokacijo določimo s tremi podatki. Ker so koordinate definirane glede na ravnino zemljine orbite, je z koordinata za Zemljo vedno enaka 0.

l longituda, ekliptična dolžina, od 0° do 360° .

b latituda, ekliptična širina od -90° do 90° .

r razdalja

V kartezične koordinate podatke pretvorimo po formulah:

$$x = r \cos b \cos l$$

$$y = r \cos b \sin l$$

$$z = r \sin b$$

Ker tiri vseh planetov ležijo v (skoraj) isti ravnini, bomo z koordinato zanemarili.