

Regresie liniară

Cursul 7

Programare și modelare probabilistă - anul III

Facultatea de Informatică, UAIC

e-mail: adrian.zalinescu@uaic.ro

web: <https://sites.google.com/view/fiicoursepmp/home>

27 Noiembrie 2023

- 1 Regresia liniară robustă
- 2 Regresia liniară ierarhică
- 3 Probleme legate de modelul liniar
 - Variabile ascunse
 - Mascarea efectului variabilelor
- 4 Interacțiuni între variabile
- 5 Dispersie variabilă

Utilizarea ipotezei că datele folosesc o distribuție Gaussiană este rezonabilă în multe cazuri.

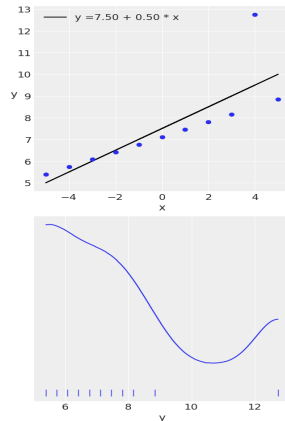
- Acest lucru nu înseamnă neapărat că datele sunt Gausseine (au fost generate conform unei distribuții normale), ci că acest lucru este o aproximare rezonabilă pentru o problemă dată.
- După cum am văzut deja în Cursul 5, presupunerea că datele sunt Gaussiene eșuează în prezența *valorilor aberante* (*outliers*).
- În acest caz, o distribuție t a lui Student este mai eficientă pentru luarea în calcul a valorilor aberante, pentru a obține o inferență mai robustă.

Exemplu: al treilea set de date din cvartetul lui Anscombe – https://en.wikipedia.org/wiki/Anscombe%27s_quartet.
Vom centra mai întâi datele:

```
ans = pd.read_csv('./data/anscombe.csv')
x_3 = ans[ans.group == 'III']['x'].values
y_3 = ans[ans.group == 'III']['y'].values
x_3 = x_3 - x_3.mean()with model:
    data_plus_one = data_generator + 1
```

Apoi, vedem cum arată acest set de date:

```
_, ax = plt.subplots(2, 1, figsize=(5, 10))
beta_c, alpha_c = stats.linregress(x_3, y_3)[:2]
ax[0].plot(x_3, (alpha_c + beta_c * x_3), 'k',
           label=f'y = {alpha_c:.2f} + {beta_c:.2f} * x')
ax[0].plot(x_3, y_3, 'Co')
ax[0].set_xlabel('x')
ax[0].set_ylabel('y', rotation=0)
ax[0].legend(loc=0)
az.plot_kde(y_3, ax=ax[1], rug=True)
ax[1].set_xlabel('y')
ax[1].set_yticks([])
plt.tight_layout()
```



Vom rescrie acum modelul, folosind o distribuție Student $t(\nu, \mu, \sigma)$ în loc de distribuția normală $\mathcal{N}(\mu, \sigma)$, unde ν este *parametrul de normalitate*.

```
with pm.Model() as model_t:
     $\alpha$  = pm.Normal(' $\alpha$ ', mu=y_3.mean(), sigma=1)
     $\beta$  = pm.Normal(' $\beta$ ', mu=0, sigma=1)
     $\epsilon$  = pm.HalfNormal(' $\epsilon$ ', 5)
     $\nu$ _ = pm.Exponential(' $\nu$ _ ', 1/29)
     $\nu$  = pm.Deterministic(' $\nu$ ',  $\nu$ _ + 1)

    y_pred = pm.StudentT('y_pred', mu= $\alpha$  +  $\beta$  * x_3,
                        sigma= $\epsilon$ , nu= $\nu$ , observed=y_3)

    idata_t = pm.sample(2000, return_inferencedata=True)
```

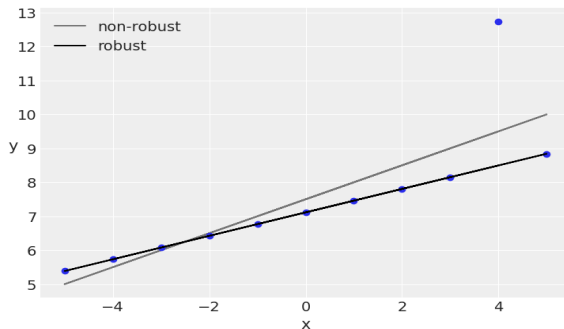
Pentru ν , distribuția a priori utilizată este $1 + \text{Exp}$ de medie 30, și nu Exp , pentru a evita valori ale lui ν prea apropiate de 0.

```

beta_c, alpha_c = stats.linregress(x_3, y_3)[:2]
plt.plot(x_3, (alpha_c + beta_c * x_3), 'k', label='non-robust', alpha=0.5)
plt.plot(x_3, y_3, 'C0o')
alpha_m = idata_t.posterior[' $\alpha$ '].mean().item()
beta_m = idata_t.posterior[' $\beta$ '].mean().item()
plt.plot(x_3, alpha_m + beta_m * x_3, c='k', label='robust')
plt.xlabel('x')
plt.ylabel('y', rotation=0)
plt.legend(loc=2)

```

O distribuție Student acordă mai puțină importanță valorilor aberante, depărtate de grosul datelor:



Rezumatul distribuției a posteriori:

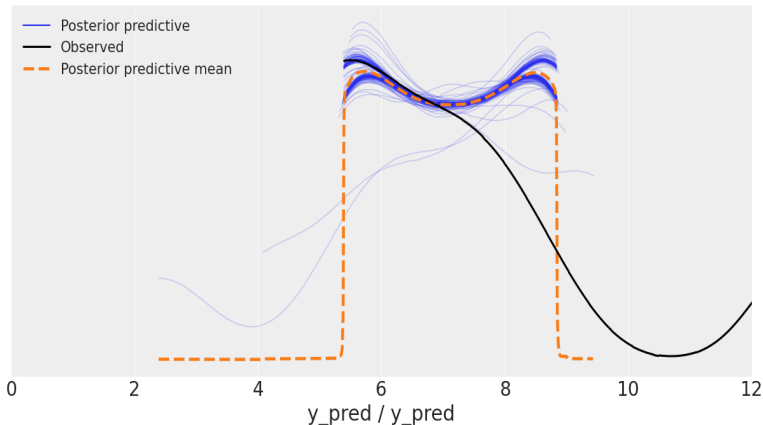
```
az.summary(idata_t, var_names=varnames)
```

	mean	sd	hdi_3%	hdi_97%	mcse_mean	mcse_sd	ess_bulk	ess_tail	r_hat
α	7.114	0.001	7.111	7.116	0.000	0.000	4721.0	4205.0	1.0
β	0.345	0.000	0.345	0.346	0.000	0.000	4554.0	3904.0	1.0
ϵ	0.003	0.002	0.001	0.006	0.000	0.000	1268.0	514.0	1.0
ν	1.211	0.201	1.000	1.570	0.003	0.002	3413.0	1591.0	1.0

Observăm că parametrii α , β și ϵ sunt foarte preciși (deviații standard foarte mici), deoarece datele erau, cu excepția valorii aberante, perfect aliniate.

Distribuția predictivă a posteriori:

```
pm.sample_posterior_predictive(idata_t, model=model_t, random_seed=2, extend_inferencedata=True)  
ax = az.plot_ppc(idata_t, num_pp_samples=200, figsize=(12, 6), mean=True)  
plt.xlim(0, 12)
```



Regresia liniară ierarhică

Vom aplica conceptul de modele ierarhice pentru regresie liniară.

- Acesta permite modelului de a considera inferențe la nivel de grup și estimări la nivel mai înalt.
- După cum am văzut în Cursul 5, acest lucru se realizează folosind *distribuții hiper a priori*.

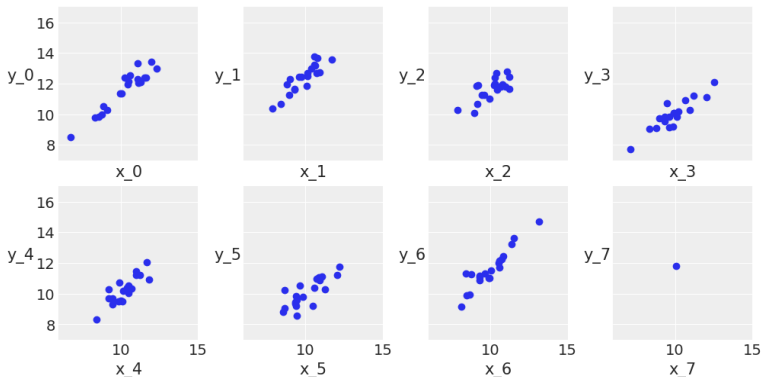
Vom crea 8 grupuri, incluzând un grup cu o singură intrare:

```
N = 20
M = 8
idx = np.repeat(range(M-1), N)
idx = np.append(idx, 7)
np.random.seed(314)

alpha_real = np.random.normal(2.5, 0.5, size=M)
beta_real = np.random.beta(6, 1, size=M)
eps_real = np.random.normal(0, 0.5, size=len(idx))

y_m = np.zeros(len(idx))
x_m = np.random.normal(10, 1, len(idx))
y_m = alpha_real[idx] + beta_real[idx] * x_m + eps_real
```

```
_, ax = plt.subplots(2, 4, figsize=(10, 5), sharex=True, sharey=True)
ax = np.ravel(ax)
j, k = 0, N
for i in range(M):
    ax[i].scatter(x_m[j:k], y_m[j:k])
    ax[i].set_xlabel(f'x_{i}')
    ax[i].set_ylabel(f'y_{i}', rotation=0, labelpad=15)
    ax[i].set_xlim(6, 15)
    ax[i].set_ylim(7, 17)
    j += N
    k += N
```



Mai întâi, centram datele: `x_centered = x_m - x_m.mean()`

Apoi, considerăm un model robust, ne-ierarhic:

```
with pm.Model() as unpooled_model:
     $\alpha_{\text{tmp}}$  = pm.Normal(' $\alpha_{\text{tmp}}$ ', mu=0, sigma=10, shape=M)
     $\beta$  = pm.Normal(' $\beta$ ', mu=0, sigma=10, shape=M)
     $\epsilon$  = pm.HalfCauchy(' $\epsilon$ ', 5)
     $\nu$  = pm.Exponential(' $\nu$ ', 1/30)

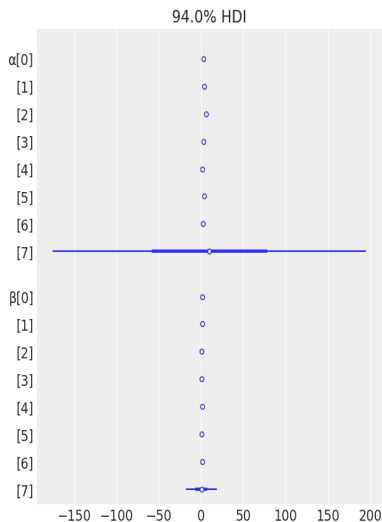
    y_pred = pm.StudentT('y_pred', mu= $\alpha_{\text{tmp}}[\text{idx}] + \beta[\text{idx}] * x_{\text{centered}}$ ,
                          sigma= $\epsilon$ , nu= $\nu$ , observed=y_m)

     $\alpha$  = pm.Deterministic(' $\alpha$ ',  $\alpha_{\text{tmp}} - \beta * x_{\text{m.mean()}}$ )

    idata_up = pm.sample(2000, target_accept=0.9, return_inferencedata=True)
```

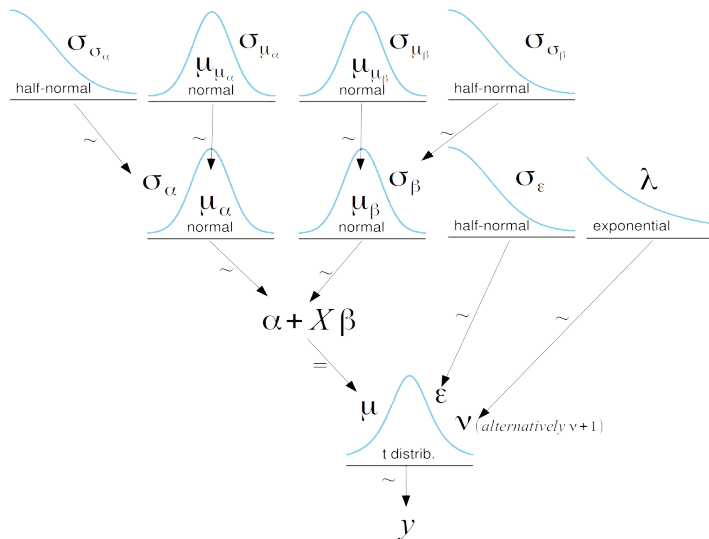
După cum vedem mai jos, estimările pentru parametrii α_7 și β_7 sunt foarte împrăștiate față de celelalte ($\alpha_i, \beta_i, 0 \leq i \leq 6$):

```
az.plot_forest(idata_up, var_names=[' $\alpha$ ', ' $\beta$ '], combined=True)
```



- Acest lucru se întâmplă deoarece nu are sens să potrivim (fit) o dreaptă numai după un punct.
- Avem nevoie de cel puțin două puncte pentru acest lucru.
- O primă posibilitate ar fi să dăm mai multe informații legate de α și β , adică să introducem distribuții tari (informative) a priori pentru acestea.
- O altă posibilitate este de a introduce un model ierarhic, ceea ce permite ca informația să fie partajată între grupuri, fapt ce conduce la fenomenul de “shrinkage” ale valorilor plauzibile ale parametrilor estimați.
- Acest lucru este foarte folositor în cazul în care avem grupuri sărace în date (cazul de mai sus este extrem: un grup cu o singură intrare).

Descrierea modelului: cu o *diagramă Krusche*:



Modelul ierarhic:

```
with pm.Model() as hierarchical_model:
    # hyper-priors
     $\alpha_{\mu\_tmp}$  = pm.Normal(' $\alpha_{\mu\_tmp}$ ', mu=0, sigma=10)
     $\alpha_{\sigma\_tmp}$  = pm.HalfNormal(' $\alpha_{\sigma\_tmp}$ ', 10)
     $\beta_{\mu}$  = pm.Normal(' $\beta_{\mu}$ ', mu=0, sigma=10)
     $\beta_{\sigma}$  = pm.HalfNormal(' $\beta_{\sigma}$ ', sigma=10)

    # priors
     $\alpha\_tmp$  = pm.Normal(' $\alpha\_tmp$ ', mu= $\alpha_{\mu\_tmp}$ , sigma= $\alpha_{\sigma\_tmp}$ , shape=M)
     $\beta$  = pm.Normal(' $\beta$ ', mu= $\beta_{\mu}$ , sigma= $\beta_{\sigma}$ , shape=M)
     $\epsilon$  = pm.HalfCauchy(' $\epsilon$ ', 5)
     $\nu$  = pm.Exponential(' $\nu$ ', 1/30)

     $y\_pred$  = pm.StudentT(' $y\_pred$ ', mu= $\alpha\_tmp[idx] + \beta[idx] * x\_centered$ ,
                        sd= $\epsilon$ , nu= $\nu$ , observed= $y\_m$ )

     $\alpha$  = pm.Deterministic(' $\alpha$ ',  $\alpha\_tmp - \beta * x\_m.mean()$ )
     $\alpha_{\mu}$  = pm.Deterministic(' $\alpha_{\mu}$ ',  $\alpha_{\mu\_tmp} - \beta_{\mu} * x\_m.mean()$ )
     $\alpha_{\sigma}$  = pm.Deterministic(' $\alpha_{sd}$ ',  $\alpha_{\sigma\_tmp} - \beta_{\mu} * x\_m.mean()$ )

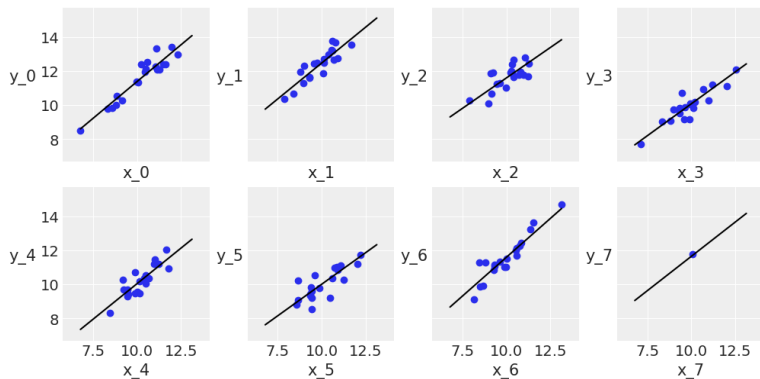
    idata_hm = pm.sample(2000, target_accept=0.99, return_inferencedata=True)
```

Vizualizarea rezultatelor:

```
_, ax = plt.subplots(2, 4, figsize=(10, 5), sharex=True, sharey=True)
ax = np.ravel(ax)
j, k = 0, N
x_range = np.linspace(x_m.min(), x_m.max(), 10)

posterior_hm = idata_hm.posterior.stack(samples={"chain", "draw"})

for i in range(M):
    ax[i].scatter(x_m[j:k], y_m[j:k])
    ax[i].set_xlabel(f'x_{i}')
    ax[i].set_ylabel(f'y_{i}', labelpad=17, rotation=0)
    alpha_m = posterior_hm[" $\alpha$ "].sel({" $\alpha_{dim_0}$ ":i}).mean().item()
    beta_m = posterior_hm[" $\beta$ "].sel({" $\beta_{dim_0}$ ":i}).mean().item()
    ax[i].plot(x_range, alpha_m + beta_m * x_range, c='k',
               label=f'y = {alpha_m:.2f} + {beta_m:.2f} * x')
plt.xlim(x_m.min()-1, x_m.max()+1)
plt.ylim(y_m.min()-1, y_m.max()+1)
j += N
k += N
```

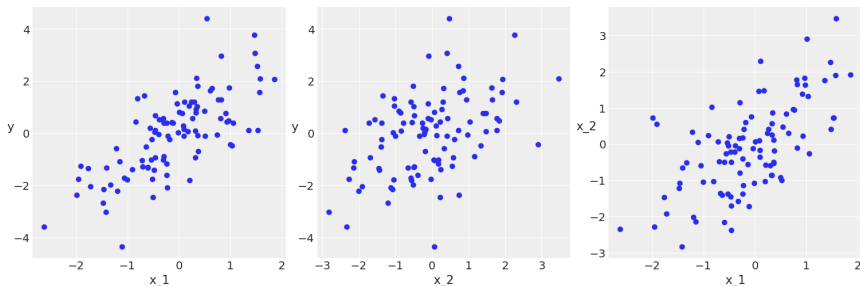
Să ne imaginăm următoarea situație:

- Fie variabila prezicătoare x și variabila prezisă y ;
- să presupunem că avem o altă variabilă z , ce influențează direct atât pe x , cât și pe y .
- De exemplu,
 - ▶ z ar putea fi “revoluția industrială”;
 - ▶ x , numărul de mori de vânt;
 - ▶ y , concentrația de CO_2 din atmosferă.
- Dacă îl omitem pe z din analiză, am putea deduce o relație liniară între x și y , ba chiar am putea prezice y în funcție de x .
- Astfel, am putea omite mecanismul care leagă aceste variabile, punând încălzirea globală pe seama reducerii numărului de mori de vânt.
- O astfel de variabilă z se numește *variabilă ascunsă* (sau *latentă*).

Următorul cod simulează o variabilă ascunsă x_1 , ce influențează variabilele x_2 și y :

```
np.random.seed(42)
N = 100
x_1 = np.random.normal(size=N)
x_2 = x_1 + np.random.normal(size=N, scale=1)
#x_2 = x_1 + np.random.normal(size=N, scale=0.01)
y = x_1 + np.random.normal(size=N)
X = np.vstack((x_1, x_2)).T
```

```
scatter_plot(X, y)
```



Vom construi trei modele:

- primul model, de regresie multiplă, cu variabilele de input x_1 și x_2 :

```
with pm.Model() as m_x1x2:
     $\alpha$  = pm.Normal('α', mu=0, sigma=10)
     $\beta_1$  = pm.Normal('β1', mu=0, sigma=10)
     $\beta_2$  = pm.Normal('β2', mu=0, sigma=10)
     $\epsilon$  = pm.HalfCauchy('ε', 5)
     $\mu$  =  $\alpha$  +  $\beta_1$  * X[:, 0] +  $\beta_2$  * X[:, 1]
    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=y)
    idata_x1x2 = pm.sample(2000, return_inferencedata=True)
```

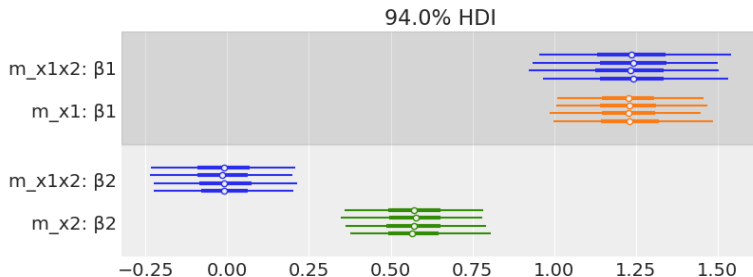
- celelalte modele, de regresie simplă, cu x_1 , respectiv x_2 :

```
with pm.Model() as m_x1:
     $\alpha$  = pm.Normal('α', mu=0, sigma=10)
     $\beta_1$  = pm.Normal('β1', mu=0, sigma=10)
     $\epsilon$  = pm.HalfCauchy('ε', 5)
     $\mu$  =  $\alpha$  +  $\beta_1$  * X[:, 0]
    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=y)
    idata_x1 = pm.sample(2000, return_inferencedata=True)

with pm.Model() as m_x2:
     $\alpha$  = pm.Normal('α', mu=0, sigma=10)
     $\beta_2$  = pm.Normal('β2', mu=0, sigma=10)
     $\epsilon$  = pm.HalfCauchy('ε', 5)
     $\mu$  =  $\alpha$  +  $\beta_2$  * X[:, 1]
    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=y)
    idata_x2 = pm.sample(2000, return_inferencedata=True)
```

Să aruncăm o privire asupra parametrilor β din aceste modele:

```
az.plot_forest([idata_x1x2, idata_x1, idata_x2],  
               model_names=['m_x1x2', 'm_x1', 'm_x2'],  
               var_names=[' $\beta_1$ ', ' $\beta_2$ '],  
               combined=False, colors='cycle', figsize=(8, 3))
```

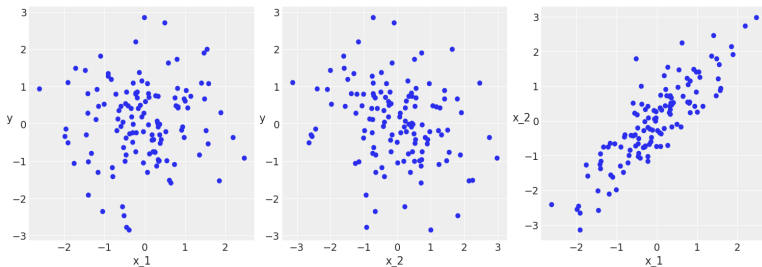


- ▶ După cum vedem, β_2 din modelul m_{x1x2} este în jurul lui 0, ceea ce indică o contribuție aproape neglijabilă a variabilei x_2 pentru a explica variabila y .
- ▶ În schimb, dacă variabila x_1 este ascunsă, β_2 devine semnificativ mai mare (peste 0.5).
- ▶ Astfel, informația în x_2 este *redundantă* dacă știm x_1 .

Mascarea efectului variabilelor

Vom crea un set de date ce exemplifică acest fenomen:

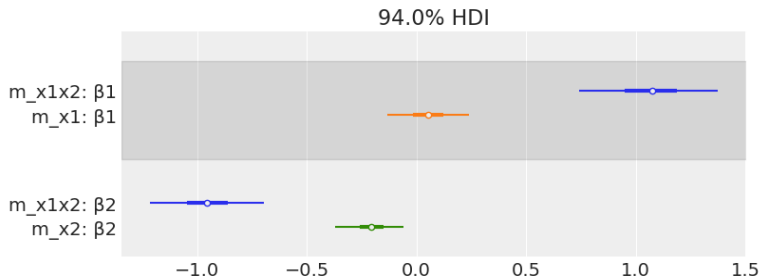
```
np.random.seed(42)
N = 126
r = 0.8
x_1 = np.random.normal(size=N)
x_2 = np.random.normal(x_1, scale=(1 - r ** 2) ** 0.5)
y = np.random.normal(x_1 - x_2)
X = np.vstack((x_1, x_2)).T
scatter_plot(X, y)
```



Astfel, x_1 și x_2 sunt pozitiv corelate, iar x_1 și x_2 sunt corelate cu y , dar în sens invers.

Construind cele trei modele, ca mai înainte, obținem următoarele rezultate:

```
az.plot_forest([idata_x1x2, idata_x1, idata_x2],  
               model_names=['m_x1x2', 'm_x1', 'm_x2'],  
               var_names=[' $\beta_1$ ', ' $\beta_2$ '],  
               combined=True, colors='cycle', figsize=(8, 3))
```



- ▶ Observăm că valorile lui β_1 și β_2 din modelul m_{x1x2} sunt în jurul lui 1, respectiv -1 , ceea ce este de așteptat.
- ▶ În schimb, în modelele simple, aceste valori sunt mai apropiate de 0, ceea ce indică un efect mai slab: are loc o anulare parțială a efectelor:
- ▶ dacă y și x_1 cresc, atunci x_2 descrește, chiar dacă x_1 și x_2 sunt pozitiv corelate.

Interacțiuni între variabile

Într-un model de regresie multiplă, se presupune (implicit) că o schimbare în x_i rezultă într-o schimbare constantă în y , dacă păstrăm fixe valorile restului variabilelor independente.

- Dar acest lucru nu este întotdeauna adevărat.
- Ce se întâmplă dacă schimbările în x_i afectează y , dar rata acestora este modulată de schimbările în x_j ?
- Exemplu privind interacțiunea între medicamente:
 - ▶ creșterea dozei medicamentului A rezultă în efecte pozitive asupra pacientului;
 - ▶ dar acest lucru este valabil în absența (sau doze mici ale) medicamentului B , în timp ce efectul lui A este negativ (sau letal) pentru doze mărite ale lui B .

Dacă dorim să captăm aceste efecte, o opțiune este să includem termeni non-aditivi în modelul nostru.

- O opțiune standard, dar nu unica, este să multiplicăm variabilele, de exemplu:

$$\mu = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2.$$

- Acest termen non-aditiv este cunoscut în statistică drept *interacțiune*.
- Putem interpreta acest termen ca un model liniar. De exemplu, putem scrie expresia de mai sus ca

$$\mu = \alpha + \underbrace{(\beta_1 + \beta_3 x_2)}_{\text{panta lui } x_1} x_1 + \beta_2 x_2$$

sau ca

$$\mu = \alpha + \beta_1 x_1 + \underbrace{(\beta_2 + \beta_3 x_1)}_{\text{panta lui } x_2} x_2.$$

Dispersie variabilă

Până acum am folosit un tipar liniar (sau polinomial) pentru a modela **media** unei variabile:

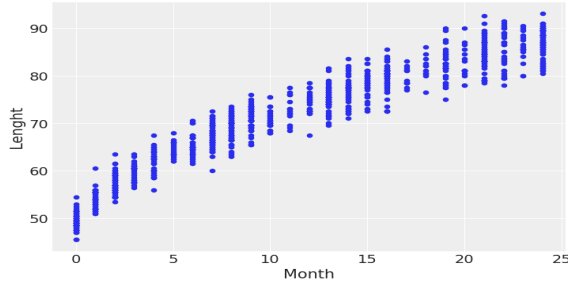
$$\mathbb{E}y = \mu = \alpha + \beta x.$$

Putem face același lucru și pentru dispersia (ori deviația standard) a variabilei dependente, atunci când nu mai putem presupune că aceasta este constantă.

În acest caz, vom considera dispersia ca o funcție (liniară) de variabila independentă.

Exemplu: OMS colectează date despre copiii din întreaga lume și concepe standarde de creștere pentru aceștia. Un exemplu de asemenea date este oferit de înălțimea copiilor de până la o anumită vârstă în funcție de vârsta lor (în luni):

```
data = pd.read_csv('./data/babies.csv')  
data.plot.scatter('Month', 'Lenght')
```



Comparativ cu modelele precedente, vom mai introduce 3 elemente suplimentare:

- ε este acum o funcție liniară de x : $\varepsilon = \gamma + \delta x$;
- astfel γ și δ sunt analoage lui α , respectiv β ;
- modelul liniar pentru μ (media) este o funcție de \sqrt{x} : $\mu = \alpha + \beta\sqrt{x}$;
- acesta este doar o modalitate de a potrivi modelul pe datele de mai sus.

- definim o variabilă partajată, `x_shared`. O vom folosi pentru a schimba valorile lui x fără a recalibra modelul.

Modelul este definit astfel:

```
with pm.Model() as model_vv:
     $\alpha$  = pm.Normal(' $\alpha$ ', sigma=10)
     $\beta$  = pm.Normal(' $\beta$ ', sigma=10)
     $\gamma$  = pm.HalfNormal(' $\gamma$ ', sigma=10)
     $\delta$  = pm.HalfNormal(' $\delta$ ', sigma=10)

    x_shared = pm.MutableData("x_shared", data.Month.values.astype(float))

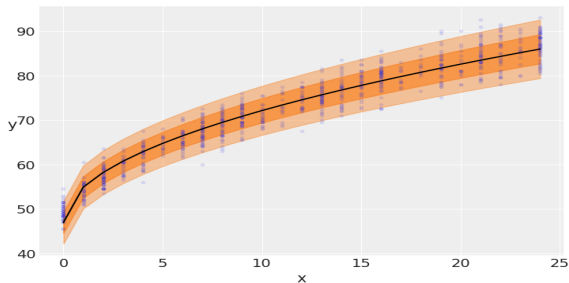
     $\mu$  = pm.Deterministic(' $\mu$ ',  $\alpha + \beta * x\_shared**0.5$ )
     $\epsilon$  = pm.Deterministic(' $\epsilon$ ',  $\gamma + \delta * x\_shared$ )

    y_pred = pm.Normal('y_pred', mu= $\mu$ , sigma= $\epsilon$ , observed=data.Lenght)

    idata_vv = pm.sample(1000, tune=1000, return_inferencedata=True)
```

Vizualizarea distribuției a posteriori:

```
plt.plot(data.Month, data.Lenght, 'C0.', alpha=0.1)
μ_m = idata_vv.posterior['μ'].mean(("chain", "draw")).values
ε_m = idata_vv.posterior['ε'].mean(("chain", "draw")).values
plt.plot(data.Month, μ_m, c='k')
plt.fill_between(data.Month,
                 μ_m + 1 * ε_m, μ_m - 1 * ε_m,
                 alpha=0.6, color='C1')
plt.fill_between(data.Month,
                 μ_m + 2 * ε_m, μ_m - 2 * ε_m,
                 alpha=0.4, color='C1')
plt.xlabel('x')
plt.ylabel('y', rotation=0)
```



- Să presupunem acum că avem un exemplu concret de copil de două săptămâni (0.5 luni) și dorim să aflăm cum se situează el în această analiză.
- Putem face acest lucru cu ajutorul funcției `sample_posterior_predictive` din PyMC3, cerând modelului distribuția variabilei $y = \text{length}$ pentru $x = 0.5$ luni.
- Problema este că această funcție returnează predicții ale lui y pentru valori observate ale lui x , iar valoarea 0.5 nu a fost observată.
- O posibilitate este de a defini o variabilă partajată (ca parte a modelului) și apoi a actualiza valoarea aceste variabile partajate chiar înainte de a genera valori din distribuția predictivă a posteriori:

```
pm.set_data({"x_shared": [0.5]}), model=model_vv)
ppc = pm.sample_posterior_predictive(idata_vv, model=model_vv)
y_ppc = ppc.posterior_predictive['y_pred'].stack(sample=("chain", "draw")).val
```

Putem acum vizualiza distribuția lui $y \mid x = 0.5$ și, adițional, cum se plasează o valoare de referință (o înălțime) relativ la aceasta:

```
ref = 53
grid, pdf = az.kde(y_ppc)
plt.plot(grid, pdf)
percentile = int(np.mean(y_ppc <= ref) * 100)
plt.fill_between(grid[grid < ref],
                 pdf[grid < ref],
                 label=f'percentile = {percentile:2d}')
plt.xlabel('length')
plt.yticks([])
plt.legend()
```

