

# Data Science Supervision 1

5. PDF of a normal distribution is  $\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} = f(x)$

$$\therefore \log(f(x)) = \text{constant} - \log(\sigma) - \frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2$$

$$\log(\text{likelihood}) = \left(\sum_{i=1}^m \log(f_A(x_i))\right) + \left(\sum_{i=1}^n \log(f_B(y_i))\right)$$

assuming independence

$$= \left(\sum_{i=1}^m -\log(\sigma) - \frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2\right) + \left(\sum_{i=1}^n -\log(\sigma) - \frac{1}{2}\left(\frac{y_i-\mu-\delta}{\sigma}\right)^2\right) + \text{constant}$$

~~$$= (m+n) \log \sigma - \frac{1}{2\sigma^2} \left( \sum_{i=1}^m x_i^2 + \sum_{i=1}^n y_i^2 \right) + \dots$$~~

$$= \text{constant} - (m+n) \log \sigma - \frac{1}{2\sigma^2} \left( \sum_{i=1}^m (x_i-\mu)^2 + \sum_{i=1}^n (y_i-\mu-\delta)^2 \right)$$

$$\frac{\partial \log(\text{likelihood})}{\partial \delta} = -\frac{1}{2\sigma^2} \sum_{i=1}^n -2(y_i - \mu - \delta)$$

$$= \frac{1}{\sigma^2} (n\bar{y} - n\mu - n\delta)$$

$$\stackrel{!}{=} 0$$

$$\therefore \bar{y} - \mu - \delta = 0$$

$$\therefore \mu + \delta = \bar{y}$$

$$\frac{\partial \log(\text{likelihood})}{\partial \mu} = -\frac{1}{2\sigma^2} \left( \sum_{i=1}^m -2(x_i - \mu) + \sum_{i=1}^n -2(y_i - \mu - \delta) \right)$$

$$= \frac{1}{\sigma^2} (m\bar{x} - m\mu + n\bar{y} - n\mu - n\delta)$$

$$= \frac{1}{\sigma^2} (m\bar{x} - m\mu + 0)$$

$$\stackrel{!}{=} 0$$

$$\therefore \bar{x} = \mu \Rightarrow \delta = \bar{y} - \bar{x}$$



$$\frac{\partial \log(\text{likelihood})}{\partial \sigma} = -\frac{m+n}{\sigma} + \frac{1}{4\sigma^3} \left( \left( \sum_{i=1}^m (x_i - \mu)^2 \right) + \left( \sum_{i=1}^n (y_i - \mu - \delta)^2 \right) \right)$$

$$= 0$$

$$\therefore 4\sigma^2(m+n) - \left( \left( \sum_{i=1}^m (x_i^2 + \mu^2 - 2\mu x_i) \right) + \left( \sum_{i=1}^n (y_i^2 + (\mu + \delta)^2 - 2(\mu + \delta)y_i) \right) \right) = 0$$

$$\therefore 4\sigma^2(m+n) - (m\bar{x}^2 + m\bar{x}^2 - 2m\bar{x} + n\bar{y}^2 + n\bar{y}^2 - 2n\bar{y}) = 0$$

$$\therefore 4\sigma^2 = \frac{m(\bar{x}^2 - \bar{x}^2) + n(\bar{y}^2 - \bar{y}^2)}{n + m}$$

$$\therefore \sigma = \frac{1}{2} \sqrt{\frac{m(\bar{x}^2 - \bar{x}^2) + n(\bar{y}^2 - \bar{y}^2)}{n+m}}$$

7. Using feature vectors:

$$\underbrace{\begin{matrix} \vec{x} \vec{1}_{x < \text{inflection}} & \vec{x} \vec{1}_{x > \text{inflection}} & \vec{1}_{x < \text{inflection}} & \vec{1}_{x > \text{inflection}} \end{matrix}}_{\text{element-wise multiplication}}$$

The linear model can be written as

$$y = m_1 \vec{x} \vec{1}_{x < \text{inflection}} + c_1 \vec{1}_{x < \text{inflection}} + m_2 \vec{x} \vec{1}_{x > \text{inflection}} + c_2 \vec{1}_{x > \text{inflection}}$$

Constrained by the equation  $m_1 \cdot \text{inflection} + c_1 = m_2 \cdot \text{inflection} + c_2$

$$\therefore c_2 = (m_1 - m_2) \cdot \text{inflection} + c_1$$

Let  $i = \text{inflection}$ .

$$\begin{aligned} \therefore y &= m_1 \vec{x} \vec{1}_{x < i} + c_1 \vec{1}_{x < i} + m_2 \vec{x} \vec{1}_{x > i} + m_1 \cdot i \vec{1}_{x > i} - m_2 \vec{1}_{x > i} + c_1 \vec{1}_{x > i} \\ &= m_1 (\vec{x} \vec{1}_{x < i} + i \vec{1}_{x > i}) + m_2 (\vec{x} \vec{1}_{x > i} - i \vec{1}_{x > i}) + c_1 (\vec{1}_{x < i} + \vec{1}_{x > i}) \end{aligned}$$



With new feature vectors

$$A = \overrightarrow{x1_{x < i}} + i \overrightarrow{1_{x \geq i}}$$

$$B = \overrightarrow{x1_{x \geq i}} + i \overrightarrow{1_{x < i}}$$

$$C = \overrightarrow{1}$$

The model can be written as

$$y = m_1 A + m_2 B + c_1 C$$

$$\text{where } c_2 = c_1 + (m_1 - m_2) i$$

8. Let  $\vec{\delta}_k = \vec{1}_{u=k}$

$$\text{temp} \triangleq \vec{a} + \beta_1 \sin(2\pi \vec{t}) + \beta_2 \cos(2\pi \vec{t}) + \epsilon_1 \vec{\delta}_{\text{decade}_1=1980s} + \epsilon_2 \vec{\delta}_{\text{decade}_1=1990s} \\ + \epsilon_3 \vec{\delta}_{\text{decade}_1=2000s} + \epsilon_4 \vec{\delta}_{\text{decade}_1=2010s} + \epsilon_5 \vec{\delta}_{\text{decade}_1=2020s}$$

Feature vectors are  $\{\vec{a}, \sin(2\pi \vec{t}), \cos(2\pi \vec{t})\} \cup \{\vec{\delta}_k \mid k \in \{\text{decade}_1=1980s, \dots, \text{decade}_1=2020s\}\}$

It can be solved using the attached code.

```
1 (a, b1, b2, e1, e2, e3, e4, e5) = linear_model.LinearRegression( fit_intercept=False).fit(
2     numpy.column.stack(
3         numpy.ones(len(temp)),
4         numpy.sin(2 * numpy.pi * t),
5         numpy.cos(2 * numpy.pi * t),
6         numpy.where(u="decade_1980s", 1, 0),
7         numpy.where(u="decade_1990s", 1, 0),
8         numpy.where(u="decade_2000s", 1, 0),
9         numpy.where(u="decade_2010s", 1, 0),
10        numpy.where(u="decade_2020s", 1, 0)
11    ),
12    temp
13 )
```



10. Let  $\vec{\gamma}_k = \vec{1}_{\text{gender} = k}$ ,

$$\vec{\delta}_k = \vec{1}_{\text{eth} = k}$$

$$\vec{1}_{\text{outcome} = \text{"fnd"}} = \epsilon_1 \vec{\delta}_{\text{Asian}} + \epsilon_2 \vec{\delta}_{\text{Black}} + \epsilon_3 \vec{\delta}_{\text{Mixed}} + \epsilon_4 \vec{\delta}_{\text{Other}} + \epsilon_5 \vec{\delta}_{\text{White}} \\ + \epsilon_6 \vec{\gamma}_{\text{female}} + \epsilon_7 \vec{\gamma}_{\text{male}}$$

The parameters are not identifiable because

$$\sum_{\text{eth}} \vec{\delta} = \cancel{\sum_{\text{eth}} \vec{\gamma}} \sum_{\text{gender}} \vec{\gamma}$$

A linearly independent model would be:

$$\vec{1}_{\text{outcome} = \text{"fnd"}} = \epsilon_1 \vec{\delta}_{\text{Asian}} + \epsilon_2 \vec{\delta}_{\text{Black}} + \epsilon_3 \vec{\delta}_{\text{Mixed}} + \epsilon_4 \vec{\delta}_{\text{Other}} \\ + \epsilon_5 \vec{\delta}_{\text{White}} + \epsilon_6 \vec{\gamma}_{\text{female}} + \epsilon_7 \left( \sum_{\text{eth}} \vec{\delta} - \vec{\gamma}_{\text{female}} \right) \\ = (\epsilon_1 + \epsilon_7) \vec{\delta}_{\text{Asian}} + (\epsilon_2 + \epsilon_7) \vec{\delta}_{\text{Black}} + (\epsilon_3 + \epsilon_7) \vec{\delta}_{\text{Mixed}} \\ + (\epsilon_4 + \epsilon_7) \vec{\delta}_{\text{Other}} + (\epsilon_5 + \epsilon_7) \vec{\delta}_{\text{White}} + (\epsilon_6 - \epsilon_7) \vec{\gamma}_{\text{female}}$$

where  $\epsilon_1$  = bias towards Asian people

$\epsilon_2$  = bias towards Black people

$\epsilon_3$  = bias towards Mixed race people

$\epsilon_4$  = bias towards people of other ethnicities

$\epsilon_5$  = bias towards white people

$\epsilon_6$  = bias towards females

$\epsilon_7$  = bias towards males