

Regresión y optimización

Reconocimiento de patrones

Gamaliel Moreno

Enero-julio 2021

UAZ-MCPI

Maestría en Ciencias del Procesamiento de la Información

1. Regresión lineal

Interpretación probabilística

2. Conclusion

Regresión lineal

- $(\mathbf{x}^{(i)}, y^{(i)})$: i-ésimo dato de entrenamiento
- $h_{\theta}(\mathbf{x}^{(i)})$: predicción de hipótesis h_{θ} para dato de entrada $\mathbf{x}^{(i)}$

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^n \theta_j x_j = \boldsymbol{\theta}^T \mathbf{x}$$

donde asumimos $x_0 = 1$. El número de características es n

- Función cuadrática de costo:

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^m h_{\theta}(\mathbf{x}^{(i)} - y^{(i)})^2$$

con m el número de datos en el conjunto de entrenamiento

- Analíticamente: $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

- ¿Por qué usamos una función cuadrática de costo?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

como m el número de datos en el conjunto de entrenamiento.

- Pudimos usar otras cosas: valor absoluto , o potencia de 4...
- Vamos a analizar el problema de regresión desde una perspectiva probabilística, como posible argumentación para esto

Regresión lineal

- Supongamos que la salida y las entradas siguen el modelo de regresión

$$y^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} + \epsilon^{(i)} \Rightarrow \epsilon^{(i)} = y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}$$

- El término de error $\epsilon^{(i)}$ captura
 - Aspectos no modelados
 - Ruido aleatorio
- Supondremos que $\epsilon^{(i)}$ son independientes e idénticamente distribuidos (i.i.d.)
- Como proceso es complejo, supondremos que $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ y por tanto la PDF de $\epsilon^{(i)}$ es

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

- Introduciendo el modelo $\epsilon^{(i)} = y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}$ derivamos

$$p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right)$$

- Nótese que θ está fuera de la condición, pues no es una variable aleatoria. La asumimos como dada.
- La distribución $y^{(i)} | \mathbf{x}^{(i)}; \theta \sim \mathcal{N}(\theta^T \mathbf{x}^{(i)}, \sigma^2)$

- Con la matriz de diseño \mathbf{X} y los parámetros de θ , la probabilidad conjunta de los datos es $p(\mathbf{y}|\mathbf{X}; \theta)$
- $p(\mathbf{y}|\mathbf{X}; \theta)$ se interpreta como función de los datos para θ constante
- Cuando queremos interpretar a $p(\mathbf{y}|\mathbf{X}; \theta)$ como función de θ la llamamos verosimilitud (likelihood)

$$L(\theta) = L(\theta; \mathbf{X}, \mathbf{y}) = p(\mathbf{y}|\mathbf{X}; \theta)$$

- Si suponemos que $\epsilon^{(i)}$ son i.i.d, entonces

$$L(\theta) = \prod_{i=1}^m p(y^{(i)}|\mathbf{x}^{(i)}; \theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right)$$

Verosimilitud logarítmica

- La verosimilitud logarítmica (log likelihood) es

$$\begin{aligned}l(\boldsymbol{\theta}) &= \ln(L(\boldsymbol{\theta})) \\&= \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right) \\&= \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right) \\&= m \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2\end{aligned}$$

- Observe que maximizar $l(\boldsymbol{\theta})$ es lo mismo que minimizar $J(\boldsymbol{\theta})$
- Con suposiciones probabilísticas: regresión de mínimos cuadrados es equivalente a estimación de máxima verosimilitud de $\boldsymbol{\theta}$ (Note irrelevancia σ)

- En ejemplo de precios de casas, tenemos varias características (features) a disposición:
 - Área habitable
 - Número de pisos
 - Número de habitaciones
- Selección de las cuáles características usar es un criterio de diseños
- Es posible introducir características artificiales para introducir no linealidad en un proceso en principio lineal:

$$\mathbf{x} = [1 \quad x_1 \quad x_1^2 \cdots x_1^n]$$

lo que permite modelar aproximaciones de Taylor de n-ésimo orden, de cualquier función no lineal.

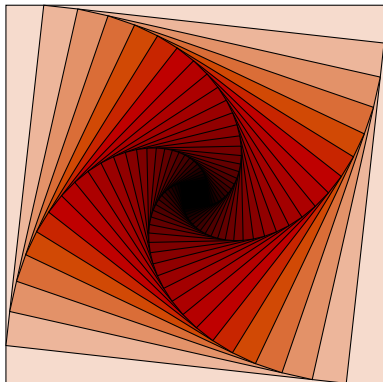


Figure 1: Rotated square from texample.net.

Table 1: Largest cities in the world (source: Wikipedia)

City	Population
Mexico City	20,116,842
Shanghai	19,210,000
Peking	15,796,450
Istanbul	14,160,467

Three different block environments are pre-defined and may be styled with an optional background color.

Default

Block content.

Default

Block content.

Alert

Block content.

Alert

Block content.

Example

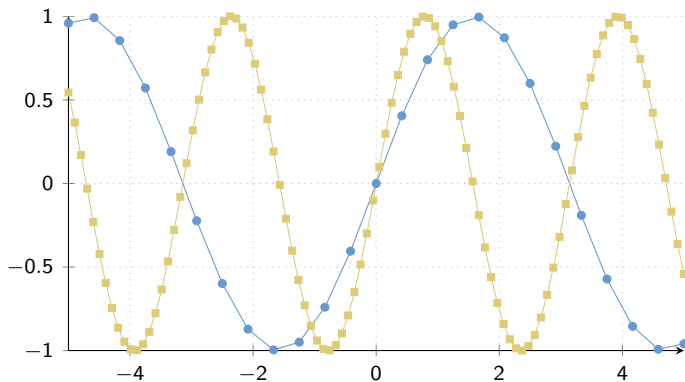
Block content.

Example

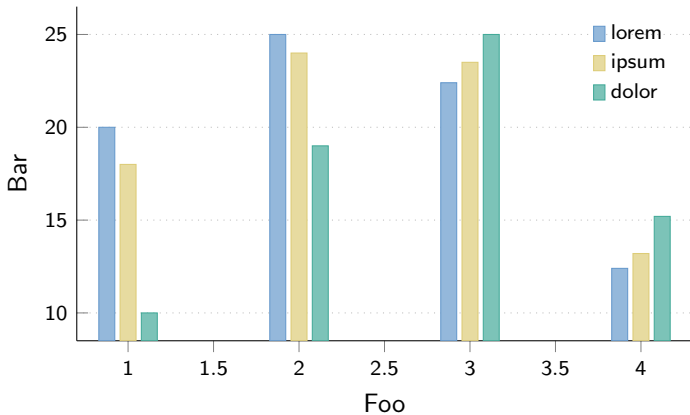
Block content.

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$

Line plots



Bar charts



Veni, Vidi, Vici

metropolis defines a custom beamer template to add a text to the footer. It can be set via

```
\setbeamertemplate{frame footer}{My custom footer}
```

Some references to showcase `[allowframebreaks]` [?, ?, ?, ?, ?]

Conclusion

Get the source of this theme and the demo presentation from

`github.com/matze/mtheme`

The theme *itself* is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.



Questions?

Sometimes, it is useful to add slides at the end of your presentation to refer to during audience questions.

The best way to do this is to include the `appendixnumberbeamer` package in your preamble and call `\appendix` before your backup slides.

metropolis will automatically turn off slide numbering and progress bars for slides in the appendix.

