

Regression

INSTRUCTOR: HONGJIE CHEN
JUNE 6TH 2022

Extend Y from Discrete to Continuous

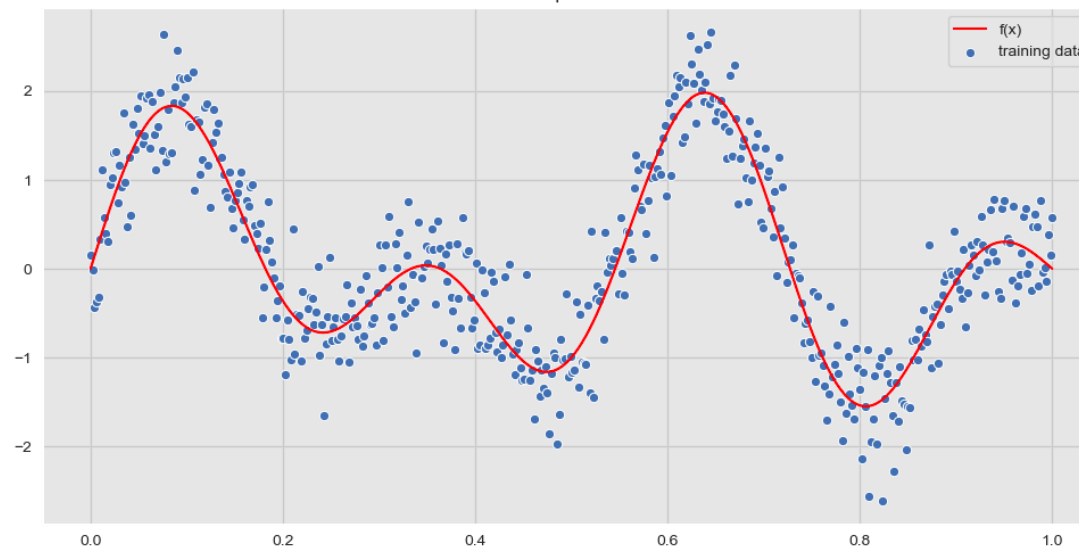
- Classification:
 $P(Y | X)$ where Y is discrete
- Regression:
 $P(Y | X)$ where Y is continuous
- For example...

Problem setting

- Learn a function $f : X \rightarrow Y, Y \in \mathcal{R}$
- Approach:
 - Choose some parameterized form for $P(Y | X, \theta)$ where θ is called a parameter vector.
 - Estimate θ via MLE or MAP

Parameterized Form for $P(Y | X, \theta)$

- Assume Y is some deterministic $f(X)$, plus random noise
 - $y = f(x) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma)$



- Therefore $p(y | x) = \mathcal{N}(f(x), \epsilon)$
- Expectation: $\mathbb{E}[Y] = f(X)$

Figure credit: [link](#)

Where we were last time

- Review the logistic regression model
- Distinguish generative models and discriminative models
 - Generative models describe how data are generated
 - Discriminative models distinguish data by boundaries
- Regression model
- HW1, HW2, HW3

Linear Regression

- $p(y | x) = \mathcal{N}(f(x), \epsilon)$
- Assume $f(x)$ is a linear function
$$p(y | x) = \mathcal{N}(w_1x + w_0, \sigma)$$
$$\mathbb{E}(y | x) = w_1x + w_0$$
- Make parameters explicit
 - $W = [w_1, w_0]$
 - $p(y | x, W) = \mathcal{N}(w_1x + w_0, \sigma)$

Training Linear Regression Model

- MCLE

$$W_{MCLE} = \operatorname{argmax}_W \prod_l P(Y^l | X^l, W)$$

$$W_{MCLE} = \operatorname{argmax}_W \sum_l \ln P(y^l | x^l, W) \quad \text{Log likelihood}$$

$$\text{Where } P(y | x, W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y - \mu)^2}{2\sigma^2}}, \mu = f(x, W)$$

Simplify

$$P(y|x, W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-f(x, W))^2}{2\sigma^2}}$$

$$W_{MCLE} = \operatorname{argmax}_W \sum_l \ln P(y^l | x^l, W)$$

$$= \operatorname{argmax}_W \sum_l \left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} + \left(-\frac{(y^l - f(x^l, W))^2}{2\sigma^2} \right) \right)$$

Constant

$$= \operatorname{argmax}_W -\frac{1}{2\sigma^2} \sum_l (y^l - f(x^l, W))^2$$

$$= \operatorname{argmin}_W \frac{1}{2\sigma^2} \sum_l (y^l - f(x^l, W))^2$$

$$= \operatorname{argmin}_W \sum_l (y^l - (w_1 x^l + w_0))^2$$

Gradient Descent

- $W_{MCLE} = \operatorname{argmin}_W \sum_l (y^l - f(x^l, W))^2$ E

- Gradient $\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$, is a vector

- Training rule: $\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} - \eta \nabla E(\vec{w})$

- View from one feature dimension $\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$

Calculate Derivative

- $E = \sum_l (y^l - (w_1 x^l + w_0))^2$

- $\frac{\partial E}{\partial w_0} =$

- $\frac{\partial E}{\partial w_1} =$

Calculate Derivative

- $E = \sum_l (y^l - (w_1 x^l + w_0))^2$
- $\frac{\partial E}{\partial w_0} = -2 \sum_l (y^l - (w_1 x^l + w_0))$
- $\frac{\partial E}{\partial w_1} = -2 \sum_l (y^l - (w_1 x^l + w_0)) x^l$

Vectorize $X = [x_1, x_2, \dots, x_n]$

$$f(x) = w_0 + \sum_{j=1}^n w_j x_j$$

$$\vec{w} = [w_0, w_1, \dots, w_n]$$

$$\frac{\partial E}{\partial w_i} = -2 \sum_l (y^l - (w_0 + \sum_{j=1}^n w_j x_j)) x_i^l$$

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} - \eta \nabla E(\vec{w})$$

$$w_i \leftarrow w_i + 2\eta \sum_l (y^l - (w_0 + \sum_{j=1}^n w_j x_j)) x_i^l$$

Gradient Descent Algorithm

- Repeat until $\Delta w < \epsilon$

- For all dimension i ,

$$w_i \leftarrow w_i + 2\eta \sum_l (y^l - (w_0 + \sum_{j=1}^n w_j x_j)) x_i^l$$

- Assume $x_0 = 1$ to incorporate w_0

MAP

- $$W_{MAP} = \operatorname{argmin}_W \left(-c \sum_i w_i^2 \right) + \sum_l (y^l - f(x^l, W))^2$$

Regularization

- Remember advantages of regularization?
- Demo: <https://lukaszkujawa.github.io/gradient-descent.html>
- Question: must f a linear function to x ?