

CO2774
Machine Learning
Sep 3, 2021

Last class:

Newton's Method, 3 way second order information for optimization
Locally weighted linear Regression :- very primitive method

$$\{x^{(i)}, y^{(i)}\}_{i=1}^m$$

$$y \in \mathbb{R} \leftarrow h_0(x) = \theta^T x$$

First Classification:
Algorithm:

$$x \in \mathbb{R}^{n+1}$$

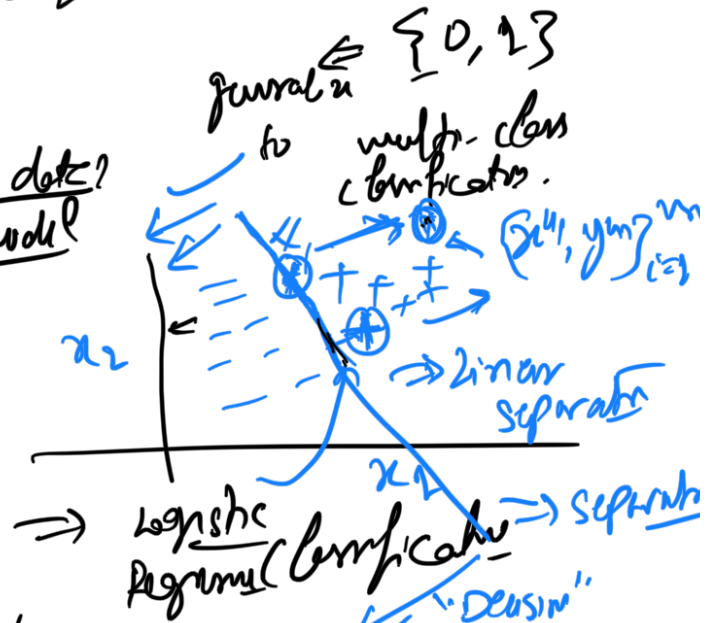
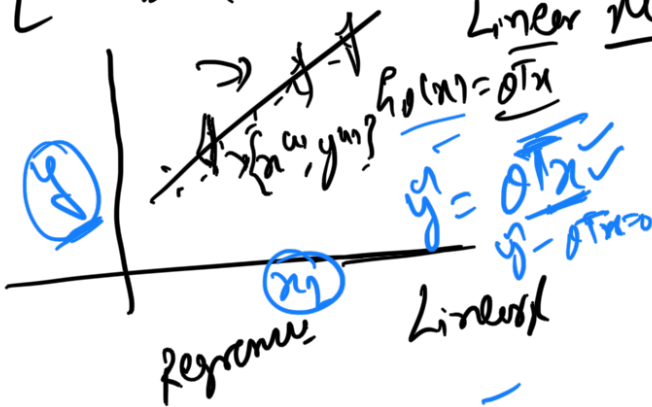
$$\{x^{(i)}, y^{(i)}\}_{i=1}^m$$

\mathbb{R} (regression)

$y \in$ Discrete Set

What kind of models can be learned over this data?

Linear Model



$$y^{(i)} | x^{(i)}, \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

$$\min_{\theta} L_2(\theta) \equiv \text{Least square}$$

Decision Boundary
 $\theta^T x = 0$
hyperplane $x \in \mathbb{R}^{n+1}$
 $\begin{cases} y = \text{white} \\ y = \text{black} \end{cases}$

$$\theta_2 x_1 + \theta_1 x_2 + \theta_0 = 0$$

$$x_2 = -\left(\frac{\theta_0 + \theta_1 x_1}{\theta_2}\right)$$

→ Classification -

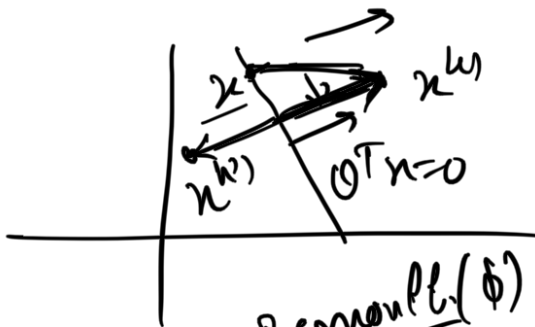
$$y^{(i)} | x^{(i)}; \theta \sim \text{Bernoulli}(\phi) \in [0, 1]$$

$\in \{0, 1\}$ ✓

→ probability of a point being +ve

ϕ to axis (signed) distance of the point from the line

$$\phi \in [0, 1] \quad \therefore \quad \theta^T x = 0$$



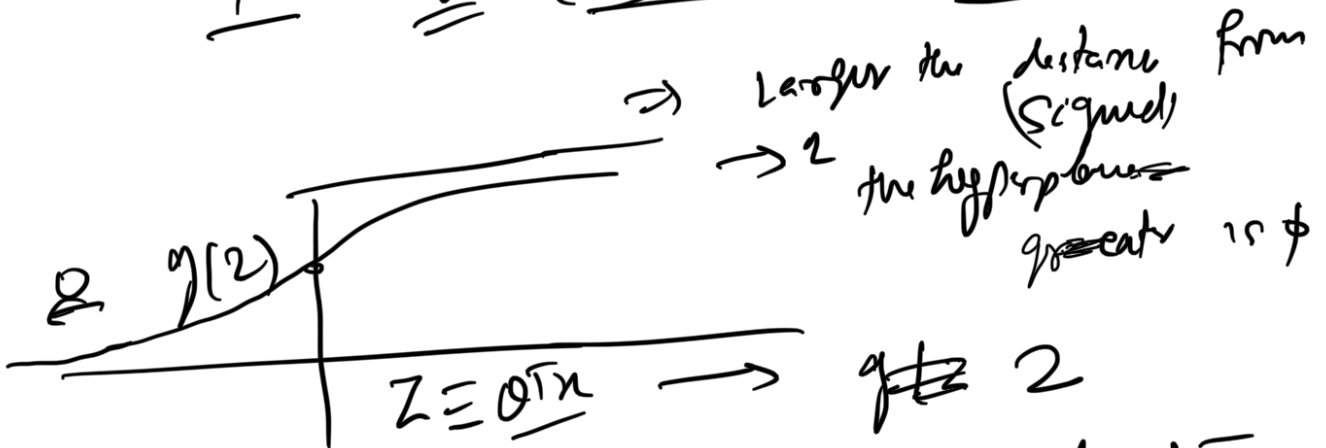
$$y^{(i)} | x^{(i)}; \theta \sim \text{Bernoulli}(\phi) \quad \theta^T x^{(i)} \quad (\text{signed})$$

$$\theta^T (x^{(i)} - x) = \text{unnormalized dist}$$

$$\theta^T x^{(i)} - \theta^T x$$

$$\theta^T x^{(i)} \equiv \text{unnormalized distance of } x^{(i)} \text{ from } \theta^T x = 0$$

$$\phi = g[\theta^T x^{(i)}] \in [0, 1]$$



S-shaped } logistic

$$\phi = g(z) = \frac{1}{1 + e^{-z}} \quad \text{sigmoid / logistic}$$

~~$y^{(i)}$~~

$$y^{(i)} | x^{(i)}; \theta \sim \text{Bernoulli}(\underbrace{g(\theta^T x^{(i)})}_{\phi^{(i)}})$$

$$P(y^{(i)}=1 | x^{(i)}; \theta) = \underbrace{g(\theta^T x^{(i)})}_{\substack{\equiv \phi^{(i)} \\ \leftarrow 1/(1+e^{-\theta^T x^{(i)}})}}$$

$$y \sim \text{Bernoulli}(\phi) \quad y \in \{0, 1\}$$

$$P(y=1) = \phi$$

$$P(y=0) = 1 - \phi$$

$$\{x^{(i)}, y^{(i)}\}_{i=1}^m$$

$$(y^{(i)} | x^{(i)}; \theta) \sim \text{Bernoulli}(\underbrace{\phi^{(i)}}_{g(\theta^T x^{(i)})})$$

→ log-likelihood of the data under the above modeling assumption.

$$P(y^{(1)} \dots y^{(m)} | x^{(1)} \dots x^{(m)}; \theta)$$

$$= \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) \equiv \text{likelihood of the data}$$

$$\Rightarrow \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

$$\sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; \theta) \quad \text{--- (1)}$$

$$= \sum_{i=1}^n$$

$$\rightarrow y^{(i)} = 1 \mid \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$

$$y \in \{0, 1\} \xrightarrow{\text{indicator fn}} y^{(i)} = 1 \text{ if } \text{Boolean expr} \text{ is true} \rightarrow y^{(i)} = 0 \text{ if false} \left\{ \begin{array}{l} 1 - \frac{1}{1 + e^{-\theta^T x^{(i)}}} \\ \frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \end{array} \right\}$$

$$= \sum_{i=1}^n \left[y^{(i)} \log p(y^{(i)} = 1 \mid x^{(i)}; \theta) + (1 - y^{(i)}) \log p(y^{(i)} = 0 \mid x^{(i)}; \theta) \right]$$

\Downarrow
 $y^{(i)} = 0$

$$\mathbb{1}\{ \text{Boolean expr} \} = \begin{cases} 1 & \text{if exp is true} \\ 0 & \text{if exp is false} \end{cases}$$

$$\begin{aligned} \text{log-likelihood} &= - \sum_{i=1}^n \left[y^{(i)} \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right] + (1 - y^{(i)}) \log \left[\frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \right] \right] \\ &= - \sum_{i=1}^n \left[y^{(i)} \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right] + (1 - y^{(i)}) \log \left[\frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \right] \right] \end{aligned}$$

What is the set of parameters that we are interested in finding?



$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(\theta)$$

$$\Rightarrow \underset{\theta}{\operatorname{argmax}} \frac{1}{2m} \sum_{i=1}^m y^{(i)} \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right] + (1 - y^{(i)}) \left[\log \left[\frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \right] \right]$$

\Rightarrow How to find θ 's?
 \hookrightarrow Gradient Descent.

$$\nabla_{\theta} \mathcal{L}(\theta) \Downarrow$$

$$\frac{1}{2m} \nabla_{\theta} \left[\sum_{i=1}^m y^{(i)} \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right] + [1 - y^{(i)}] \log \left[\frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \right] \right] \begin{matrix} A \\ B \end{matrix}$$

$$= \frac{1}{2m} \nabla_{\theta} \left[\sum_{i=1}^m \frac{y^{(i)} \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right]} + [1 - y^{(i)}] \log \left[\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right] + [1 - y^{(i)}] \log (e^{-\theta^T x^{(i)}}) \right]$$

$$= \frac{1}{2m} \nabla_0 \left[\sum_{i=1}^m \log \left[\frac{1}{1+e^{-\theta^T x^{(i)}}} \right] + [1-y^{(i)}] (-\theta^T x^{(i)}) \right]$$

$$= \frac{1}{2m} \nabla_0 \left[\sum_{i=1}^m -\log [1+e^{-\theta^T x^{(i)}}] + (1-y^{(i)}) (-\theta^T x^{(i)}) \right]$$

$$= \frac{1}{2m} \sum_{i=1}^m -\nabla_0 \log [1+e^{-\theta^T x^{(i)}}] + (1-y^{(i)}) \nabla_0 (-\theta^T x^{(i)})$$

$$= \frac{1}{2m} \sum_{i=1}^m - \left[\frac{1 \cdot (e^{-\theta^T x^{(i)}}) (-x^{(i)})}{(1+e^{-\theta^T x^{(i)}})} + (1-y^{(i)}) (-x^{(i)}) \right]$$

$$= \frac{1}{2m} \sum_{i=1}^m \left[1-y^{(i)} - \left[\frac{e^{-\theta^T x^{(i)}}}{1+e^{-\theta^T x^{(i)}}} \right] \right] (-x^{(i)})$$

$$= \frac{1}{2m} \sum_{i=1}^m \left[\cancel{e^{-\theta^T x^{(i)}}} \left[1 - \frac{e^{-\theta^T x^{(i)}}}{1+e^{-\theta^T x^{(i)}}} \right] + y^{(i)} \right] (-x^{(i)})$$

$$\nabla_0 L(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[y^{(i)} - \frac{1}{1+e^{-\theta^T x^{(i)}}} \right] (-x^{(i)})$$

11

12