

Lecture - 1

Machine Learning

gives computer the ability to learn without being explicitly programmed.

Supervised Learning

for given dataset with input x and label y ,
then learn mapping from x to y .

Regression problem - value y is continuous

Classification problem - y takes on discrete number of variables

Support vector machine \Rightarrow infinite dimension

Deep learning

a neural network with 3 or more layers

Unsupervised learning

dataset with ~~no~~ labels (i.e. x with no y given)

* Use x and find interesting things about it.

Reinforcement learning

Let it find the way it wants and if it's correct output, then it's fine but if it shows wrong output, we correct it.

Lecture-2

Linear Regression

Training set

↓
Learning algorithm

↓
 h (hypothesis)

function
takes input
& tries to tell output.

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

↑
parameter
feature 1

↑
parameter
feature 2

$$h(x) = \sum_{j=0}^n \theta_j x_j$$

where $x_0 = 1$
(dummy feature)

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

θ = parameter

m = training example
/ no of rows
in table

x = 'input' / feature

y = output / target
variable

(x, y) = training example

(x^i, y^i) = i^{th} training
example

$$h_{\theta}(x) = h(x)$$

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

minimize $J(\theta)$

Gradient Descent

Start with some θ (say $\theta = \vec{0}$)
keep changing θ to reduce $J(\theta)$.

$$\theta_j \stackrel{\text{assignment}}{:=} \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (j = 0, 1, 2, \dots)$$

learning rate

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y)$$

$$= (h_{\theta}(x) - y) \cdot x_j$$

Repeat until convergence

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x_j^i$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

(for $j = 0, 1, \dots, n$)

Stochastic gradient descent

repeat

$$\left\{ \begin{array}{l} \text{for } i = 1 \text{ to } m \\ \theta_j \leftarrow \theta_j - \alpha (h_{\theta}(x^i) - y^i) x_j^i \\ \end{array} \right\}$$

for every j

* Normal eqⁿ

set $\theta_{\text{old}} = \theta$ and in one step with a few matrix multiplications you end up with the optimal value of θ that land you right at the global optimum (in 1 step)

$$\theta = \begin{pmatrix} - \\ - \\ - \end{pmatrix}$$

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

$\theta \in \mathbb{R}^{n+1}$ (3-D)

$$Q = \frac{1}{2} \sum_{i=1}^n (h(x^i) - y^i)^2$$

$$XQ = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \vdots \\ (x^m)^T \end{bmatrix} \begin{bmatrix} Q_0 \\ Q_1 \\ Q_2 \end{bmatrix}$$

$$= \begin{bmatrix} (x^1)^T Q \\ (x^2)^T Q \\ \vdots \\ (x^m)^T Q \end{bmatrix}$$

$$= \begin{bmatrix} h_Q(x^1) \\ \vdots \\ h_Q(x^m) \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$J(Q) = \frac{1}{2} (XQ - y)^T (XQ - y)$$

$$x\theta - y = \begin{bmatrix} h_{\theta}(x^1) - y^1 \\ \vdots \\ h_{\theta}(x^m) - y^m \end{bmatrix}$$

$$(x\theta - y)^T (x\theta - y)$$

$$\left[\because z^T z = \sum z^2 \right]$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T x^T - y^T) (x\theta - y)$$

$$= \frac{1}{2} \nabla_{\theta} [\theta^T x^T x\theta - \theta^T x^T y - y^T x\theta + y^T y]$$

$$= \frac{1}{2} [x^T x\theta + x^T x\theta - x^T y - x^T y]$$

$$= x^T x\theta - x^T y \stackrel{\text{set}}{=} \vec{0}$$

$$x^T x\theta = x^T y \quad \left[\text{"Normal eq"} \right]$$

$$\theta = (x^T x)^{-1} x^T y$$

~~00000~~

Lecture - 3

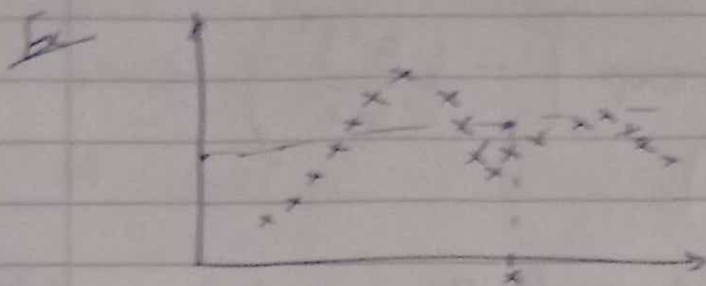
Locally weighted Regression

"Parametric" learning algorithm

Fit fixed set of parameters (θ) to data

"Non-parametric" learning algorithm

Amount of data/parameter you need to keep grows (linearly) with size of data



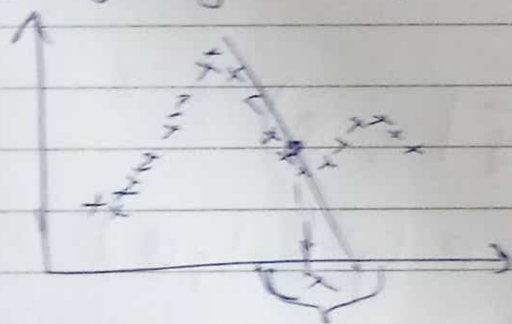
To evaluate h at certain x :

LR: Fit θ to minimize

$$\frac{1}{2} \sum_i (y^i - \theta^T x^i)^2$$

Return $\theta^T x$

Locally weight regression



(τ = bandwidth)

Fit θ to minimize

$$\sum_{i=1} w^i (y^i - \theta^T x^i)^2$$

where w^i is a "weighting" fn

$$w^i = \exp \left(-\frac{(x^i - x)^2}{2\tau^2} \right)$$

If $|x^i - x|$ is small, $w^i \approx 1$

If $|x^i - x|$ is large, $w^i \approx 0$

why τ^2 ?

* Why least squares?

Assume $y^i = \mathcal{Q}^T x^i + \varepsilon^i$
 ε^i includes unmodelled effects, random noise

$\varepsilon^i \sim \mathcal{N}(0, \sigma^2)$
 \uparrow twiddle \Rightarrow distributed as

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

i.i.d. \Rightarrow independently & Identically Distributed
 \uparrow given

$$P(y^i | x^i; \mathcal{Q}) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^i - \mathcal{Q}^T x^i)^2}{2\sigma^2}\right)$$

\uparrow mean
 \uparrow variance

i.e. $y^i | x^i; \mathcal{Q} \sim \mathcal{N}(\mathcal{Q}^T x^i, \sigma^2)$

~~likelihood of \mathcal{Q}~~

"likelihood of \mathcal{Q} "

$$\begin{aligned} L(\mathcal{Q}) &= P(\vec{y} | x; \mathcal{Q}) \\ &= \prod_{i=1}^m P(y^i | x^i; \mathcal{Q}) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^i - \mathcal{Q}^T x^i)^2}{2\sigma^2}\right) \end{aligned}$$

Likelihood OR Probability

- If you view this as a fⁿ of the parameters holding the data fixed, then we call that the Likelihood.
- So if ~~to~~ we think as training set as a fixed thing and then varying parameters θ , then we use term "likelihood".
- Whereas if we view the parameters θ as fixed and maybe varying the data, we are gonna say probability.
- ~~*~~ Likelihood of parameter
Probability of data

"log likelihood"

$$\ell(\theta) = \log L(\theta)$$

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

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

$$= m \log \frac{1}{\sqrt{2\pi} \sigma} + \sum_{i=1}^m -\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}$$

MLE: maximum likelihood estimation

Choose θ to maximize $L(\theta)$

i.e. Choose θ to minimize $\frac{1}{2} \sum_{i=1}^n (y^i - \theta^T x^i)^2 = J(\theta)$

Classification

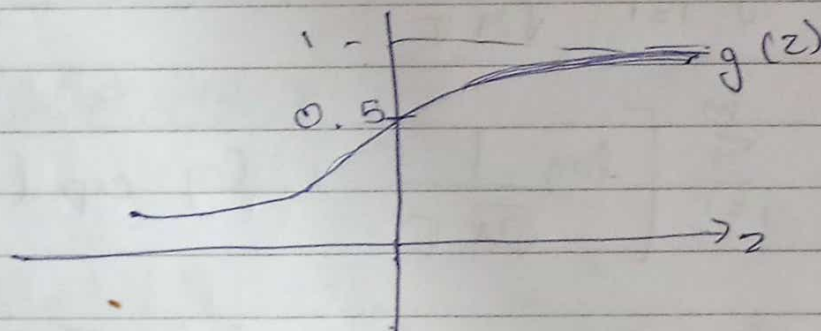
Logistic regression

want $h_{\theta}(x) \in [0, 1]$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

"sigmoid" or "logistic" function



$$P(y=1 | x; \theta) = h_{\theta}(x)$$

$$P(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

$$y \in \{0, 1\}$$

$$P(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$$

$$L(\theta) = P(\tilde{y} | x; \theta)$$

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

$$= \prod_{i=1}^m h_{\theta}(x^i)^{y^i} (1 - h_{\theta}(x^i))^{1-y^i}$$

$$\ell(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^m y^i \log h_{\theta}(x^i) + (1-y^i) \log(1-h_{\theta}(x^i))$$

Choose θ to maximize $\ell(\theta)$

Batch gradient ~~ascent~~ ascent:

$$\theta := \theta_j + \alpha \frac{d}{d\theta_j} \ell(\theta)$$

$$\theta_j = \theta_j + \alpha \underbrace{\sum_{i=1}^m (y^i - h_{\theta}(x^i)) x_j^i}_{\frac{d}{d\theta_j} \ell(\theta)}$$