

Artificial Intelligence/Machine Learning/Deep Learning: 'Bridging the Skills Gap'

Optional: Normal Equation

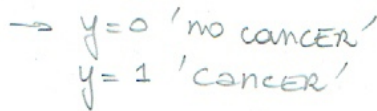
There is an analytical/closed form solution to a linear regression problem with Mean Squared Error as cost function → no gradient descent needed!

So why don't we use this all the time?

- Because it is heavy taxation on the memory → requires d^2 memory → avoid when dimension d of feature vector is large
- $X^T X$ requires lots of memory if dimension of X is large
- Inverting $X^T X$ is also complex
-

①

y_i IS A CATEGORICAL VARIABLE



values between 0 and 1 \rightarrow NOT Good!

6 win tones
↓
4 Hs, 2 Ts

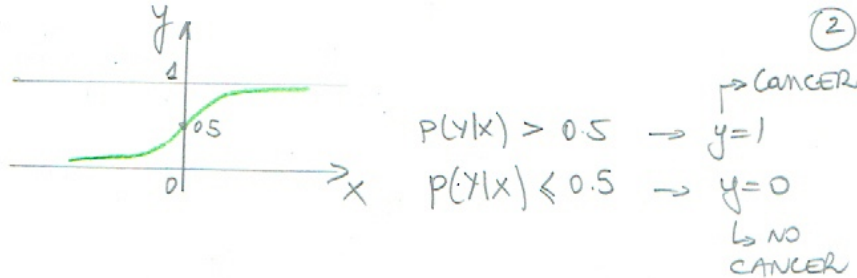
$$\ln \rightarrow \ln\left(\frac{p}{1-p}\right) = w_0 + w_1 x_1 + \dots + w_d x_d \rightarrow \text{linear function of features } x_i$$

$$\rightarrow \frac{p}{1-p} = e^{+w^T x} \rightarrow p = (e^{+w^T x}) (1-p)$$

$$\rightarrow e^{+wT_X} = p(1 + e^{+wT_X}) = e$$

$$\Rightarrow P = \frac{e^{+WTx}}{1+e^{+WTx}} \cdot \frac{e^{-WTx}}{e^{-WTx}} \rightarrow P = \frac{1}{1+e^{-WTx}}$$

probability of success



Naive Bayes $P(X|Y) \sim P(X|Y) \cdot P(Y)$

you get $P(X|Y)$ by estimating $P(X|Y)$ and $P(Y)$
 \hookrightarrow PRIORS
 likelihoods

\updownarrow
Logistics Regression

$$P(Y|X) = \frac{1}{1 + e^{-W^T X}} = \sigma(W^T X_i)$$

sigmoid

Gradient $D = \{(x_1, y_1), \dots, (x_d, y_d)\}$ m samples!
 $y_i \sim \text{Bernoulli}$

MLE: $\text{ARGMAX}_W P(D|W) \rightarrow$ Find W_s that maximizes the likelihood of seeing the DATA D

$$\begin{aligned}
 P(D|W) &= \prod_{i=1}^m P(y^{(i)} | x^{(i)}, W) \\
 &= \prod_{i=1}^m \alpha^{(i) y^{(i)}} (1 - \alpha^{(i)})^{1 - y^{(i)}} \quad y_i \in \{0, 1\} \\
 &\quad \uparrow \\
 &\quad \text{pdf Bernoulli}
 \end{aligned}$$

SET $\alpha^{(i)} = \sigma(W^T x^{(i)})$

\hookrightarrow We will not be able to solve for W because non-linearity of σ

→ Newton's method $\begin{matrix} \nabla \\ \text{H} \end{matrix}$ → 2nd order optimization (3)

- NO α
- Very Fast convergence

$$w_1 = w_0 + \frac{c'(w_0)}{c''(w_0)}$$

→ H^{-1} → can be challenge to calculate

newton if $d < 1000$

SET

$$\ell(w) \stackrel{\downarrow}{=} -\log p(D|w) \rightarrow \text{WE WANT TO MINIMIZE!}$$

- log-likelihood

$$= -\sum_{i=1}^m y_i \log \alpha_i + (1-y_i) \log (1-\alpha_i)$$

INTERMEZZO: $\frac{\partial}{\partial w_j} \log \alpha_i = + \frac{x_j e^{-w^T x}}{1 + e^{-w^T x}} = \boxed{x_j (1-\alpha)}$

$$\log \alpha = \log \sigma(w^T x) = \log \frac{1}{1 + e^{-w^T x}} = 0 - \log (1 + e^{-w^T x})$$

$$\log (1-\alpha) = -w^T x - \log (1 + e^{-w^T x})$$

$$\frac{\partial}{\partial w_j} \log (1-\alpha) = -x_j + x_j \alpha = \boxed{-\alpha x_j}$$

$$\rightarrow d(w) = - \sum_{i=1}^m y_i \log x_i + (1-y_i) \log (1-x_i)$$

$$\frac{\partial d(w)}{\partial w_j} = - \sum_{i=1}^m y_i x_{ij} (1-x_i) - (1-y_i) x_{ij} (-x_i)$$

$$= - \sum_{i=1}^m y_i x_{ij} - y_i x_{ij} x_i - x_{ij} x_i + y_i x_{ij} x_i$$

DOT product of column j of X with (x-y)

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} \end{pmatrix}$$

design matrix

1st row (x_1)^T

d-features $\rightarrow x_{i1}, \dots, x_{id}$

m-samples

i-th sample!

$$\rightarrow (x-y)^T X \rightarrow ((x-y)^T X)^T = \boxed{X^T (x-y)}$$

$\nabla_w d$

next: Hessian

$$H_w \rightarrow \frac{\partial^2}{\partial w_j \partial w_k} d_w$$

jk entry of H

$$\frac{\partial \log x}{\partial x} = \frac{1}{x}$$

$$H_w \Rightarrow \sum_{i=1}^m x_{ij} \frac{\partial}{\partial w_k} x_i (1)$$

$$\rightarrow \frac{\partial x}{\partial w_k} = x \frac{\partial \log x}{\partial x} = x \frac{1}{x} = 1$$

$$(1) = \sum_{i=1}^m \underbrace{x_{ij} x_{ik}}_{\text{row } i} x_i (1-x_i)$$

$$= z_j^T B z_k$$

DIAGONAL

WITH $z_j = (x_{1j}, \dots, x_{mj})$

$z_k = (x_{1k}, \dots, x_{mk})$

WITH $B = \begin{pmatrix} x_1(1-x_1) & & \\ & \ddots & \\ & & x_n(1-x_n) \end{pmatrix}$

$$x_j = (x_{j1}, \dots, x_{jd})^T \rightarrow \text{row } (j)$$

$$z_j = x_{1j}, \dots, x_{mj} \rightarrow \text{column } (j)$$

$$X^T B X \quad \text{to make it a row}$$

$$\rightarrow \nabla_w^2 L = X^T B X$$

positive semi
definite!

L is convex

$$\alpha_i = \nabla(w^T x_i)$$

\hookrightarrow always > 0 and < 1

$$0 < \alpha_i (1 - \alpha_i) < 1$$

Newton: \rightarrow iterative Reweighted least Squares!

\hookrightarrow Fast!

$$w_{t+1} = w_t - H^{-1} \nabla \quad (\text{see page 3})$$

$$\rightarrow w_{t+1} = w_t - (X^T B A)^{-1} \cdot X^T (x - y)$$

Assume is invertible

$$= w_t = (X^T B X)^{-1} X^T B (A w_t - B^{-1} (x - y))$$

$$= w_t = (X^T B X)^{-1} X^T B z_t$$