

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

Optional: Normal Equation

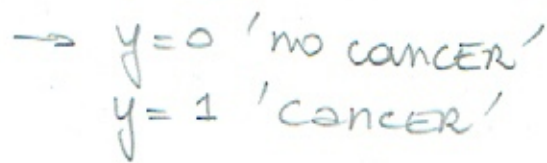
There is an analytical/closed form solution to a linear regression problem with Mean Squared Error as cost function \rightarrow no gradient descent needed!

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

- Because it is heavy taxation on the memory \rightarrow requires d^2 memory \rightarrow 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 < 0
 $> 1 \rightarrow$ NOT GOOD!

values between 0 and 1 \rightarrow NOT Good!

6 coin tosses
↓
4 Hs, 2 Ts

$$\rightarrow \frac{4}{2} \rightarrow 2 \text{ to } 1$$

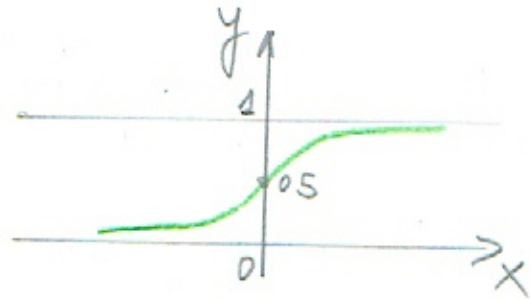
Bernoulli $B \sim p$ \leftarrow Probability of success $\rightarrow \frac{4/6}{2/6} \rightarrow \frac{4}{2} \rightarrow 2$
 vs Probability of failure.

$$\ln \rightarrow \ln\left(\frac{p}{1-p}\right) = w_0 + w_1x_1 + \dots + w_dx_d \rightarrow \text{linear function of features } x_i$$

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

$$\rightarrow e^{+wT_x} = P(1 + e^{+wT_x})$$

$$\rightarrow P = \frac{e^{+wTx}}{1+e^{+wTx}} \cdot \frac{e^{-wTx}}{e^{-wTx}} \rightarrow P = \frac{1}{1+e^{-wTx}} \quad \text{of success}$$



(2)

$$P(Y|X) > 0.5 \rightarrow y=1 \rightarrow \text{CANCER}$$

$$P(Y|X) \leq 0.5 \rightarrow y=0 \rightarrow \text{NO CANCER}$$

Naive Bayes

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

you get $P(Y|X)$ by estimating $P(X|Y)$ and $P(Y)$

likelihoods

Priors

↕

Logistics Regression

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

Sigmoid

Gradient

$$D = \{(x_1, y_1), \dots, (x_d, y_d)\} \quad m \text{ samples!}$$

$y_i \sim \text{Bernoulli}$

MLE:

$\text{Argmax}_W P(D|W) \rightarrow$ Find W s that maximizes the likelihood of seeing the DATA D

$$P(D|W) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}, W)$$

$$= \prod_{i=1}^m \alpha^{(i) y^{(i)}} (1 - \alpha^{(i)})^{1 - y^{(i)}} \quad y_i \in \{0, 1\}$$

↑
pdf Bernoulli

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

\rightarrow 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

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

- log-likelihood

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

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

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

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

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

(4)

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

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

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

$$= \sum_{i=1}^m (x_{ij} - y_i) x_{ij} \alpha_i$$

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

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

1st row (x₁)^T
design matrix

d-features → x_{i1}, ..., x_{id}
m-samples
i-th sample!

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

$\nabla_w d$

next: Hessian

$$\mathcal{H}_w \rightarrow \frac{\partial^2}{\partial w_j \partial w_k} d_w$$

jk entry
of \mathcal{H}

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

$$\mathcal{H}_w \Rightarrow \sum_{i=1}^m x_{ij} \frac{\partial}{\partial w_k} \alpha_i \quad (1)$$

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

$$(1) = \sum_{i=1}^m x_{ij} x_{ik} \alpha_i (1-\alpha_i)$$

$$\frac{\partial}{\partial w_k} \alpha_i = x_{ik} \alpha_i (1-\alpha_i)$$

↓
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 $B = \begin{pmatrix} x_1(1-x_1) & & 0 \\ & \ddots & \\ 0 & & x_n(1-x_n) \end{pmatrix}$

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

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

(5)

$$X^T B X$$

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

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

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

Normal Equation

Multi-variate regression

$$y_p = w^T x \quad (1)$$

MEAN SQUARED ERROR: $\|e\|^2 = (y_p - y)^2 \quad (2)$

(1) and (2) $\rightarrow (w^T x - y)^T (w^T x - y) \quad \|e\|^2 = \bar{e} \cdot \bar{e}^T = \bar{e}^T \bar{e}$

$$\rightarrow ((w^T x)^T - y^T) (w^T x - y)$$

$$= (w^T x)^T (w^T x) - y^T (w^T x) - y (w^T x)^T + \cancel{y^T y}$$

$$= \underbrace{x^T w w^T x}_{\|w\|^2} - \underbrace{y^T (w^T x) - (w^T x)^T y}_{-2(w^T x)^T y} \quad \begin{array}{l} \rightarrow \text{drop as no} \\ \text{function of } w. \end{array}$$

m samples
d Features.

$$\rightarrow \nabla_w C = 2x^T x w - 2x^T y$$

SET
 $\nabla_w C \downarrow = 0$

$$\rightarrow \cancel{x^T x} w - \cancel{x^T y} = 0$$

$$\rightarrow x^T x w = x^T y$$

$$\rightarrow \boxed{w = (x^T x)^{-1} x^T y}$$

pseudo-inverse
of X

\rightarrow sometimes $(x^T x)^{-1}$ is not defined

ex. $(x^T x)$ is SINGULAR (features are related)