

Recap.

1. D-Trees

2. N-Bayes.

Linear Models

(51)

- Naive Bayes "under certain conditions" is a linear model.
 - (-) So, you learn a set of weights.
- Perceptrons, "first ANN", a linear classifier model.

Linear Models :- ~~Naive Bayes~~ Basic model:-

$$w_1 \bar{x}_1 + \dots + w_k \bar{x}_k - w_0 \geq 0. \quad (\text{Perceptron})$$

Question :- How to learn the weights?

already seen two ways of doing it.

- Can be posed as an optimization problem.

$$w_{\text{opt}} = \arg \min_{\bar{w}, w_0} J(\bar{w}, T, w_0) \quad \begin{array}{l} \xrightarrow{\textcircled{1}} \text{detour to} \\ \xrightarrow{\textcircled{2}} \# \text{mistakes} \end{array}$$

Choosing a different 'F' gives a different \bar{w}_{opt}

NB :- $P(x|y), P(y) \longleftrightarrow P(y|x)$.

Logistic Regression :- Come up with a model of $P(y|x)$.

$$\text{Prob}(y=0|\bar{x}) = \frac{1}{1 + \exp(\bar{w} \cdot \bar{x} + w_0)}.$$

and,

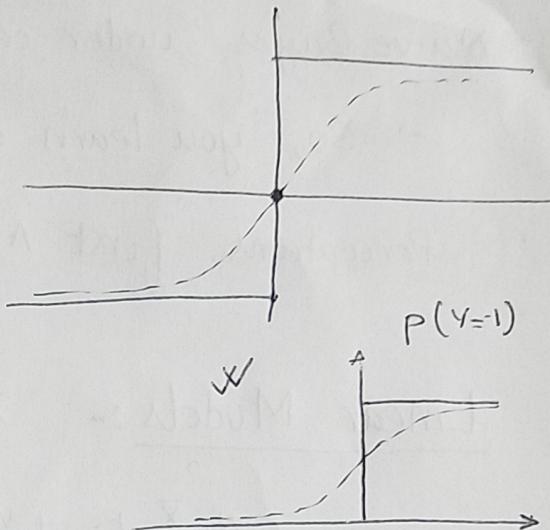
$$P(Y=1 | \bar{x}) = \frac{\exp(w_0 + \bar{w} \cdot \bar{x})}{1 + \exp(w_0 + \bar{w} \cdot \bar{x})}$$

Classification rule :-

$$P(Y=1 | \bar{x}) \geq P(Y=0 | \bar{x})$$

$$\Leftrightarrow \exp(\bar{w}^* \cdot \bar{x}^*) \geq 1$$

$$\Leftrightarrow \boxed{\bar{w} \cdot \bar{x} + w_0 \geq 0}$$



how to train our logistic regression function?

$$\bar{w} \leftarrow \arg \max_{\bar{w}} \prod_{i=1}^n P(f(\bar{x}(i)) | \bar{x}(i), \bar{w}^*)$$

For simplicity, we will take 'f' to be $\{+1, 0\}$ instead

of $\{+1, -1\}$. $\boxed{(1+Y)(1-Y) = 0}$

$$Y \in \{+1, 0\}, \text{ then, } \boxed{Y \cdot (1-Y) = 0}$$

(Ans)

$$P(f(\bar{x}) | \bar{x}, \bar{w}) = \left(P(f(\bar{x}) = 1 | \bar{x}, \bar{w}) \right)^{f(\bar{x})} \cdot \left(P(f(\bar{x}) = 0 | \bar{x}, \bar{w}) \right)^{1-f(\bar{x})}$$

\Rightarrow Log probability of the training examples given \bar{w} :-

$$\begin{aligned}
 \ell(\bar{w}) &= \sum_{i=1}^n f(\bar{x}(i)) \cdot \log \text{Prob}(f(\bar{x}(i))=1 \mid \bar{w}, \bar{x}(i)) + (1 - f(\bar{x}(i))) \cdot \log \text{Prob}(f(\bar{x}(i))=0 \mid \bar{x}, \bar{w}) \\
 &= \sum_{i=1}^n f(\bar{x}(i)) \log \frac{P(f(\bar{x}(i))=1 \mid \bar{w}, \bar{x}(i))}{P(f(\bar{x}(i))=0 \mid \bar{w}, \bar{x}(i))} + \log P(f(\bar{x}(i))=0 \mid \bar{x}, \bar{w}) \\
 &= \sum_{i=1}^n f(\bar{x}(i)) \cdot \left(w_0 + \sum_{j=1}^K w_j \bar{x}_j(i) \right) - \log \left(1 + \exp(w_0 + \sum_{j=1}^K w_j \bar{x}_j(i)) \right)
 \end{aligned}$$

NO CLOSED FORM APPROACH/SOL

BUT :- unique solution Concave function.

$\bar{w} \rightarrow \underbrace{\log(1 + \exp(\bar{x} \cdot \bar{w}))}_{\text{convex function of } \bar{w}}$.

Gradient ascent methods :-

$$\begin{aligned}
 \frac{\partial \ell(\bar{w})}{\partial w_i} &= \sum_{i=1}^n f(\bar{x}(i)) \cdot \bar{x}_j(i) - \left[\frac{\exp(\bar{w}^* \cdot \bar{x}^*)}{P(f(\bar{x}(i))=1 \mid \bar{x}, \bar{w})} \right] \cdot \bar{x}_j(i) \\
 &= \sum_{i=1}^n \bar{x}_j(i) \cdot (f(\bar{x}(i)) - P(f(\bar{x}(i))=1 \mid \bar{x}, \bar{w}))
 \end{aligned}$$

$$\bar{w}^{(l)} = \bar{w}^{(l-1)} + \eta \cdot \nabla_{\bar{w}} \ell(\bar{w})$$

$$\boxed{\bar{w}_j^{(l)} = \bar{w}_j^{(l-1)} + \eta \cdot \frac{\partial \ell(\bar{w})}{\partial w_i}}$$

$\eta \rightarrow$ large (overfitting)
 $\eta \rightarrow$ small (underfitting).

If prob close to '1' not much effect.

$\rightarrow \Theta$ A natural measure :-

Find the \bar{w} that minimizes the # of errors.

$$\bar{w}_{\text{opt}}^* = \arg \min_{\bar{w}^*} \left\{ \# \text{examples s.t. } f(\bar{x}(i)) \neq \text{sign}(\bar{w}^* \cdot \bar{x}(i)) \right\}$$

hard to compute.

$$= " " \sum_{i=1}^n \mathbb{I}\{f(\bar{x}(i)) \neq \text{sign}(\bar{w}^* \cdot \bar{x}(i))\}$$

OPTIMIZATION,

minimize (OR maximize) some $J(w_0, w_1, \dots, w_k, T)$
 ↓
 parameters to tune.

Simple example :- Regression (Linear).

- Output $f(\bar{x})$ is a real value.



- Suppose only one feature, $w_1, (\bar{x}_1) = x$.

$$\text{Area} = x \xrightarrow{\text{given}} w_0 + w_1 x.$$

$$(x_1, f_1), (x_2, f_2), \dots, (x_n, f_n).$$

Fit a line :- $y = h(x) = w_0 + w_1 x$.

Least squared error :- $\sum_{i=1}^n (f_i - (w_0 + w_1 x))^2$

Find w_0, w_1 to minimize the error.

- Convex functions :- unique minima.

$$\frac{\partial}{\partial w_0} J(w_0, w_1, T) = 0$$

$$\Rightarrow 0 = \sum_{i=1}^n 2(f_i - (w_0 + w_1 x_i)) \cdot (-1)$$

$$\Rightarrow \sum_{i=1}^n f_i = \sum_{i=1}^n (w_0 + w_1 x_i)$$

$$\Rightarrow n \cdot w_0 + w_1 \cdot (\sum x_i) = (\sum f_i)$$

$$\frac{\partial}{\partial w_1} J(w_0, w_1, T) = 0 \Rightarrow 0 = \sum_{i=1}^n 2(f_i - (w_0 + w_1 x_i)) \cdot (\sum x_i)$$

$$\Rightarrow w_0 \cdot (\sum_{i=1}^n x_i) + w_1 \cdot (\sum x_i^2) = (\sum f_i x_i)$$

$$w_1 = \frac{n \cdot (\sum f_i x_i) - (\sum f_i)(\sum x_i)}{(\sum x_i^2) - (\sum x_i)^2}$$

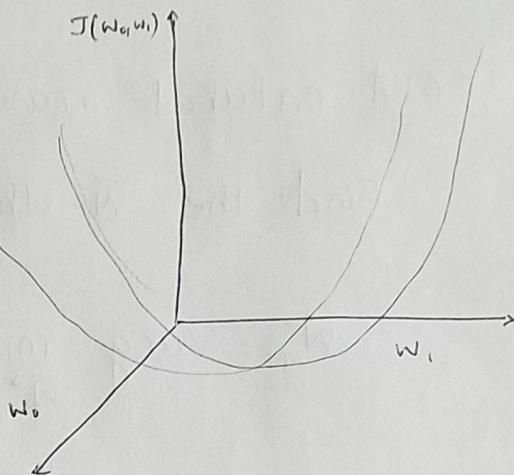
$$w_0 = \frac{\sum f_i - w_1 (\sum x_i)}{n}$$

Multivariate regression :- More features than just '1':

e.g., area, # bedrooms, walk-score.

$$f = w_0 + w_1 \bar{x}_1 + w_2 \bar{x}_2 + \dots + w_k \bar{x}_k = \bar{w} \cdot \bar{x}^* = \bar{w}^T \bar{x}^*$$

$$J(\bar{w}^*, T) = \sum_{i=1}^n (f_i - \bar{w}^* \cdot \bar{x}(i))^2 \rightarrow \bar{w}^* = (\bar{x}^T \bar{x})^{-1} \bar{x}^T \bar{y}$$



$$\min_{\bar{w}^*} \left\| F - \bar{X} \cdot \bar{w}^* \right\|_2^2$$

||

$$F \rightarrow n \times 1$$

$$\bar{X} \rightarrow n \times (k+1)$$

$$\bar{w}^* \rightarrow (k+1) \times 1$$

$$F^T = \bar{w}^{*T} (\bar{X}^T \bar{X}) \cdot w$$

$$(F^T - \bar{w}^{*T} \cdot \bar{X}^T) \cdot (F \cdot \bar{X} \cdot w) = F^T F - (\bar{w}^{*T} \cdot \bar{X}^T \cdot F) - (F^T \cdot \bar{X} \cdot w) + \bar{w}^{*T} (\bar{X}^T \bar{X}) \bar{w}^*$$

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

are the same

Take gradient and set to zero.

$$(\bar{X}^T \cdot X) \cdot \bar{w}_{opt} = \bar{X}^T \cdot F \rightarrow \text{Normal equations.}$$

- Check they hold for $k=1$ by explicitly plugging them.

$$\bar{w}_{opt}^* = (\bar{X}^T \bar{X})^{-1} \cdot \bar{X}^T \cdot F$$

MLE interpretation :-

Suppose data is generated as:-

$$f(\bar{X}) = w_0 + w_1 \bar{X}_1 + \dots + \bar{w}_k \cdot \bar{X}_k + \underbrace{\epsilon}_{N(0, \sigma^2)}$$

gaussian.

$$P(f_1, f_2, \dots, f_n | \bar{X}, \bar{w}^*, \sigma^2)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(f_i - \bar{X}_i^T \bar{w}^*)^2}{2\sigma^2}\right)$$

$$\therefore \hat{f}_i = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \cdot \exp\left(-\frac{1}{2\sigma^2} \cdot \left\| F - \bar{X} \cdot \bar{w}^* \right\|_2^2\right)$$

Gradient descent now