THEORY

OF

LOGISTIC REGRESSION

This page intentionally left blank

Gaussian Naive Bayes

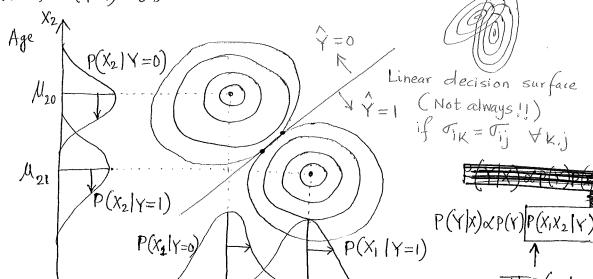
- Consider boolean $Y \in \{0,1\}$ e.g. Y = Plays Basketball?
- Consider continious X_i e.g. $X_1 = \text{Height}$, $X_2 = \text{Age}$.
- Consider P(Y=0) = P(Y=1) = 0.5

Non-spherical Gaussians.

TTP(X:1x=k)

 $\sim \mathcal{N}(\mathcal{M}_{i,k}, \mathcal{T}_i)$

3



M 10

$$Y \leftarrow \underset{y \in \{0,1\}}{\operatorname{arg max}} P(Y=y) \prod_{i} P(X_i | Y=y)$$

Parametric Form of P(YIX)

$$P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)} = \frac{P(B) = \sum P(B|Aj)P(Aj)}{if A \in \{0,1\}}$$

$$P(A|B) = \frac{P(A|B)P(Aj)}{if A \in \{0,1\}}$$

Dividing both Nr and Dr by Nr gives:

$$P(Y=|X) = \frac{1}{1 + \frac{P(Y=0) P(X|Y=0)}{P(Y=1) P(X|Y=1)}}$$
or,
$$P(Y=|X) = \frac{1}{1 + \exp\left(\ln \frac{P(Y=0) P(X|Y=0)}{P(Y=1) P(X|Y=1)}\right)}$$

$$X = \exp(\ln x)$$

 $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

if $A \in \{0,1\}$ $P(A|B) = \frac{P(B|A)P(A)}{P(A)P(B|A)+P(A')P(B|A')}$

Height

$$\frac{\sum \left(\frac{\mu_{10} - \mu_{11}}{\sigma_{12}^{2}} X_{1} + \frac{\mu_{11}^{2} - \mu_{10}^{2}}{2\sigma_{12}^{2}}\right)}{1 + \exp\left(\ln \frac{P(Y=0)}{P(Y=1)} + \ln \frac{P(X|Y=0)}{P(X|Y=1)}\right)}$$

$$\frac{1}{1 + \exp\left(\ln \frac{1-\pi}{17} + \sum \ln \frac{P(X|Y=0)}{P(X|Y=1)}\right)}$$

$$C. Indep.$$

$$\begin{array}{c} \left[\sum\limits_{i} l_{n} \frac{P[X_{i}|\gamma=0)}{P[X_{i}|\gamma=1)} \right] = \sum\limits_{i} \frac{1}{l_{n}} \frac{\frac{1}{\sqrt{2\pi\sigma_{i}^{2}}} \exp\left(\frac{-\left(X_{i}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}}\right)}{\frac{1}{\sqrt{2\pi\sigma_{i}^{2}}} \exp\left(\frac{-\left(X_{i}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}}\right)} \end{array} \\ = \sum\limits_{i} l_{n} \exp\left(\frac{\left(X_{i}-\mu_{i_{1}}\right)^{2} - \left(X_{i}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}} \right) = \sum\limits_{i} \left(\frac{\left(X_{i}-\mu_{i_{1}}\right)^{2} - \left(X_{i}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{\left(X_{i}-\mu_{i_{1}}\right)^{2} - \left(X_{i}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}} \right) = \sum\limits_{i} \left(\frac{\left(X_{i}^{2}-\mu_{i_{1}}\right)^{2} - \left(X_{i}^{2}-\mu_{i_{0}}\right)^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{2X_{i}\left(\mu_{i_{0}}-\mu_{i_{1}}\right) + \mu_{i_{1}}^{2} - \mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) = \sum\limits_{i} \left(\frac{\mu_{i_{0}}-\mu_{i_{1}}}{\sigma_{i}^{2}} \times \frac{\mu_{i_{1}}^{2}-\mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{2X_{i}\left(\mu_{i_{0}}-\mu_{i_{1}}\right) + \mu_{i_{1}}^{2} - \mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) = \sum\limits_{i} \left(\frac{\mu_{i_{0}}-\mu_{i_{1}}}{\sigma_{i}^{2}} \times \frac{\mu_{i_{1}}^{2}-\mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{\mu_{i_{0}}-\mu_{i_{1}}}{\sigma_{i}^{2}} \times \frac{\mu_{i_{0}}^{2}-\mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{\mu_{i_{0}}-\mu_{i_{0}}}{2\sigma_{i}^{2}} \times \frac{\mu_{i_{0}}-\mu_{i_{0}}^{2}}{2\sigma_{i}^{2}} \right) \\ = \sum\limits_{i} \left(\frac{\mu_{i}}{\sigma_{i}} \times \frac{\mu_{i}}{\sigma_{i}} \times \frac{\mu_{i}}{\sigma_{i}} \times \frac{\mu_{i}}{\sigma_{i}} \times \frac{\mu_{i}}{\sigma_{i}} \right) \\ = \sum\limits_{i} \left(\frac{\mu_{i}}{\sigma_{i}} \times \frac$$

 $\begin{aligned} & = \sum_{\ell} \{Y^{\ell} \ln \frac{p(Y^{\ell}=1|X^{\ell},W)}{p(Y^{\ell}=0|X^{\ell},W)} + (1-Y^{\ell}) \ln P(Y^{\ell}=0|X^{\ell},W) \} \\ & = \sum_{\ell} Y^{\ell} \ln \frac{p(Y^{\ell}=1|X^{\ell},W)}{p(Y^{\ell}=0|X^{\ell},W)} + \ln P(Y^{\ell}=0|X^{\ell},W) \\ & = \sum_{\ell} Y^{\ell} \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) - \ln \left(1 + \exp \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) \right) \\ & = \sum_{\ell} Y^{\ell} \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) - \ln \left(1 + \exp \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) \right) \\ & = \sum_{\ell} Y^{\ell} \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) - \ln \left(1 + \exp \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) \right) \\ & = \sum_{\ell} Y^{\ell} \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) - \ln \left(1 + \exp \left(\omega_{0} + \sum_{i=1}^{n} \omega_{i} \times_{i}^{i} \right) \right) \end{aligned}$

Logistic Regression More Generally

- Y is not boolean, but still discrete-valued
- y E & y,, ... yR3 : learn R-1 sets of weights

for
$$K < R$$
 $P(Y=YK|X) = \frac{\exp(\omega_{K}\phi + \frac{n}{j=1}\omega_{Ki}X_{i})}{1 + \sum_{j=1}^{R-1} \exp(\omega_{j}\phi + \sum_{i=1}^{n}\omega_{ji}X_{i})}$

for
$$K = R$$
 $P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(\omega_j \phi + \sum_{i=1}^{n} \omega_{ji} X_i)}$

Training Logistic Regression: MCLE

- we have L training examples $\{(x', Y'), ... (x', Y')\}$
- maximum likelihood estimate for parameters W

$$W_{MLE} = \underset{W}{\text{arg max}} P(\langle x', y' \rangle \dots \langle x^{L}, Y^{L} \rangle | W)$$

$$= \underset{W}{\text{arg max}} TT P(\langle x^{L}, y^{L} | W)$$

- maximum conditional likelihood estimate

Steps for training logistic regression:

- Choose parameters $W = \langle w_0, w_1 \dots w_n \rangle$ to maximize cond. likelihood of traing data where $P(Y=0|X) = \frac{1}{1+\exp(w_0 + \sum w_i X_i)}$; $P(Y=1|X) = \frac{\exp(w_0 + \sum w_i X_i)}{1+\exp(w_0 + \sum w_i X_i)}$
- Training data $D = \{(x', x'), \dots (x^L, Y^L)\}$
- Data conditional likelihood = T > (Y | X | W)

$$P(Y=\phi|X,W) = \frac{1}{1+\exp(\omega_0 + \sum_{i} \omega_i X_i)}$$

$$P(Y=1|X,W) = \frac{\exp(\omega_0 + \sum_{i} \omega_i X_i)}{1+\exp(\omega_0 + \sum_{i} \omega_i X_i)}$$

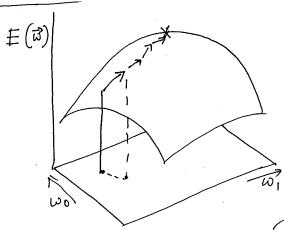
$$l(w) = ln \operatorname{TT} P(Y^{l}|X^{l}, W)$$

$$= \sum_{l} Y^{l}(\omega_{0} + \sum_{i}^{n} \omega_{i} X_{i}^{l}) - ln(1 + exp(\omega_{0} + \sum_{i}^{n} \omega_{i} X_{i}^{l}))$$

Good News! ((W) is a concave function of W

News: no closed-form solution to maximize ((W)

Ascent: Gradient



Gradient
$$\nabla E(\vec{w}) = \begin{bmatrix} \frac{\partial E}{\partial \omega_0}, \frac{\partial E}{\partial \omega_0}, \dots \frac{\partial E}{\partial \omega_n} \end{bmatrix}$$

Training Rule:
$$\overrightarrow{w}^{(i+1)} \in \overrightarrow{w}^{(i)} + \eta \nabla E(\overrightarrow{w})$$

Step Size

i.e.
$$\Delta \omega_i = \eta \frac{\partial E}{\partial \omega_i}$$

Batch gradient:

use value of the objective function $E_{p}(\vec{w})$ over the entire training data D

Do Repeat

Deat 1) Calculate
$$\nabla E_{D}(\vec{w}) = \left[\frac{\partial E_{D}(\vec{w})}{\partial \omega_{D}}, \dots \frac{\partial E_{D}(\vec{w})}{\partial \omega_{D}} \right]$$

2) Update parameters
$$\overrightarrow{w} \leftarrow \overrightarrow{w} + \eta \nabla E_{D}(\overrightarrow{w})$$

Stocknostic gradient use value of the objective function Ed (W) ove single example d ED

Do Repeat

2) Calculate the gradient just for
$$d$$

$$\nabla E_{d}(\vec{w}) = \begin{bmatrix} \partial E_{d}(\vec{w}) \\ \partial w_{0} \end{bmatrix}, \dots \partial E_{d}(\vec{w}) \end{bmatrix}$$

N.B: - Stochastic approximates batch arbitrarily closely as
$$\eta \rightarrow \beta$$

- Stochastic is much faster When D is very large
- An intermediate approach can be to choose das a subset of D

Maximize Conditional Log Likelihood: Gradient Ascent

$$\frac{12e^{-i\omega_{i}}}{l(w)} = \ln \prod_{i} P(Y^{l}|X^{l}, w)$$

$$= \sum_{i} Y^{l}(w_{o} + \sum_{i} w_{i}, x_{i}^{l}) - \ln \left(1 + \exp(w_{o} + \sum_{i} \omega_{i} x_{i}^{l})\right)$$

$$\frac{\partial \ell(w)}{\partial w_i} = \sum_{\ell} x_i^{\ell} \left(x^{\ell} - \stackrel{\wedge}{P} \left(x^{\ell} = 1 \mid x^{\ell}, w \right) \right)$$

missessigned prob. mass given current sets of weights by feature value

Gradient ascent algorithm: iterate until change < E

∀i, repeat

$$w_i \leftarrow w_i + \eta \sum_{k} x_i^{\ell} \left(Y^{\ell} - \hat{P}(Y^{\ell-1}|X^{\ell},W) \right)$$
assume $X_{p} = 1$ for w_0

We used M(c) LE. Can we use MAP?

- One common approach is to define priors on W
- Helps avoid very large weights and overfitting
- MAP estimale

AP estimate

$$W \leftarrow arg max ln P(W) TT P(Y^{l}|X^{l},W)$$
 $W \leftarrow arg max ln P(W) TT P(Y^{l}|X^{l},W)$

- Let's assume Gaussian prior WN (\$, T)

Gradient ascent agorithm: iterate until change < E

$$w_i \leftarrow w_i - \left[\eta \lambda w_i \right] + \eta \sum_{k} x_i^{k} \left(Y^{\ell} - \hat{P} \left(Y^{k-1} | X, w \right) \right)$$

regularization

(this is pushing the weights to be zero)