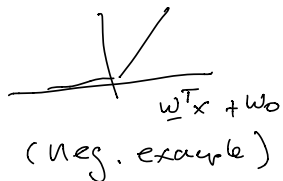


Recall $\tilde{y}_n = \begin{cases} +1 & \text{if } \underline{w}^T \underline{x}_n + w_0 > 0 \\ -1 & \text{otherwise} \end{cases}$

train via hinge loss 

But, what if we'd like a probability that an example is positive?

Approach 1: Discriminative

MLE

$$\arg \max_{\underline{w}} \prod_n p(y_n | \underline{x}_n, \underline{w})$$

conditional likelihood

⊗ simple

→ Logistic regression

Approach 2: Generative

MLE

$$\arg \max_{\underline{w}} \prod_n p(\underline{x}_n, y_n | \underline{w})$$

joint likelihood

⊗ flexible ⊗ missing labels
⊗ add knowledge

→ multivariate Gaussian

→ naive Bayes

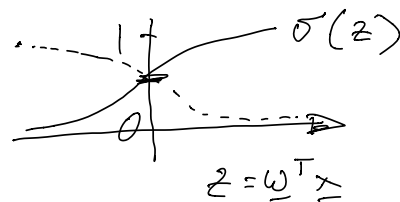
Note ① convenient to adopt $\{0, 1\}$

$$\textcircled{2} \hat{y} = \begin{cases} 1 & \text{if } p(y=1 | \underline{x}) > p(y=0 | \underline{x}) \\ 0 & \text{otherwise} \end{cases}$$

(1) Discriminative (Logistic regression)

Model $p(y|x)$ through a sigmoid (logistic)

$$p(y=1|x) = \frac{1}{1 + \exp(-\underline{w}^T \underline{x})}$$



Write shorthand $h = \underline{w}^T \underline{x}$.

Note $p(y=1|x) = \frac{1}{1+e^{-h}}$ $p(y=0|x) = \frac{1}{1+e^h}$

Since $y \in \{0, 1\}$, can write

$$p(y|x) = p(y=1|x)^y p(y=0|x)^{(1-y)}$$

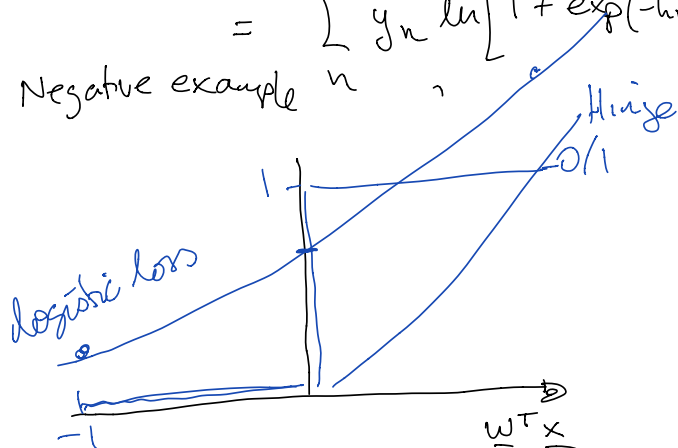
MLE Loss is negative log likelihood

$$L_D(\underline{w}) = - \sum_n \ln [p(y_n | \underline{x}_n)]$$

$$= - \sum_n y_n \ln [\sigma(h_n)] - \sum_n (1-y_n) \ln [\sigma(-h_n)]$$

$$= \sum_n y_n \ln [1 + \exp(-h_n)] + \sum_n (1-y_n) \ln [1 + \exp(h_n)]$$

Negative example



loss on a -ve ex

Differentiable & convex

Prefers to make better decisions on correct predictions

Stochastic gradient descent

Positive example

Hint: If $\hat{y}_n \neq y_n$:

$$\underline{w}_{t+1} \leftarrow \underline{w}_t + \eta \underline{x}_n$$

Negative example

If $\hat{y}_n \neq y_n$

$$\underline{w}_{t+1} \leftarrow \underline{w}_t - \eta \underline{x}_n$$

Logistic: $\underline{w}_{t+1} \leftarrow \underline{w}_t + \eta \underline{x}_n p(y_n=0 | \underline{x}_n)$

$$\underline{w}_{t+1} \leftarrow \underline{w}_t + \eta \underline{x}_n p(y_n=1 | \underline{x}_n)$$

Algorithms come from loss functions

Loss + SGD \rightarrow Algorithm

$$\underline{w}_{t+1} \leftarrow \underline{w}_t - \eta \frac{\partial L(\underline{w})}{\partial \underline{w}}$$

SGD

[2] Generative approach

Model $p(x, y)$ directly

- (1) Decide on a generative process
- (2) Parametrize model
- (3) Minimize ~~the~~ negative log. likelihood

Use the chain (product) rule:

$$p(x, y) = p(x|y) p(y)$$

class
conditional

continuous closure

Gaussian Naive Bayes

class prior

Bernoulli

(y) → (x)

[MLE] $\arg \max_{\underline{w}} \prod_n p(x_n, y_n) = \arg \max_{\underline{w}} \prod_n p(x_n|y_n) p(y_n)$

$$= \arg \min_{\underline{w}} - \underbrace{\sum_n \ln[p(x_n|y_n)]}_{\text{solve for parts of } \underline{w}^* \text{ corresponding to class conditional}} - \underbrace{\sum_n \ln[p(y_n)]}_{\text{solve for parts of } \underline{w}^* \text{ corresponding to class prior}}$$

Aside Classify via Bayes rule

$$p(y=1 | x) \propto p(y=1) p(x|y=1)$$

Estimate
Class prior

Model as Bernoulli r.v. with
prob θ

$$p(y) = \theta^y (1-\theta)^{(1-y)} \quad \text{Set } \frac{dL(\theta)}{d\theta} = 0$$

~~MLE~~ MLE: $\theta^* = \frac{\sum y_n}{N}$

$$L(\theta) = - \sum_n y_n \log(\theta) - \sum_n (1-y_n) \log(1-\theta)$$

∴

solve analytically

Class conditional (Continuous)

$p(\underline{x}|y)$ \underline{x} continuous

$$\underline{x}|y=0 \sim N(\underline{\mu}_0, \Sigma_0) \quad \underline{x}|y=1 \sim N(\underline{\mu}_1, \Sigma_1)$$

$$\underline{w} = \{ \underline{\mu}_0, \Sigma_0, \underline{\mu}_1, \Sigma_1 \} \quad (+ \theta)$$

Intuition: MLE separates by class

For class $y=0$, use

$$\{ \underline{x} : (\underline{x}_n, y_n) \in D, y_n = 0 \}$$

\hookrightarrow ~~est~~ estimate $\underline{\mu}_0, \Sigma_0$

For class $y=1$, use

$$\{ \underline{x} : (\underline{x}_n, y_n) \in D, y_n = 1 \}$$

\hookrightarrow estimate $\underline{\mu}_1, \Sigma_1$

Example

$$\text{MLE } \underline{\mu}_0 = \frac{\sum_{n: y=0} \underline{x}_n}{N_0} \quad \begin{array}{l} \# \text{ examples} \\ \leftarrow \text{ with } y_n=0 \end{array}$$

\hookrightarrow Understand decision boundaries!

(Note) Same covariance $\Sigma_1 = \Sigma_2$ then
 linear decision boundaries,
 otherwise get quadratic boundaries.

Classify $\hat{y} = \begin{cases} 1 & \text{if } p(y=1)p(x|y=1) > \\ 0 & \text{otherwise} \end{cases}$
 $p(y=0)p(x|y=0)$

Boundary:
 $S = \{x : \frac{p(x|y=1)}{p(x|y=0)} = \text{constant}\}$

(1) Equivalently, decision boundary ~~is~~ S :
 $\ln[p(x|y=1)] - \ln[p(x|y=0)] = \text{constant}$

Plug in
 $N(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$

& algebra

(see section notes)

Class conditional (discrete) "Naive Bayes"

Now x_d takes on one of $\{1, \dots, J\}$ discrete values. e.g., hair color

Key assumption (Naive Bayes)

Each dimension of \underline{x} is independent, conditioned on the class

$$p(\underline{x} | y=1) = \prod_d p(x_d | y=1)$$

For feature x_d , model as a categorical

distribution



Notation

$$\pi_{kdj} \geq 0$$

Probability in class k of feature d taking on value j

$$\sum_j \pi_{kdj} = 1$$

for all k ,
all d

Write x_d as a "one-hot" vector

$$x_d = [1 \ 0 \ 0]^T \text{ for blond hair}$$

Can write

$$P(\underline{x} | y=1) = \prod_d \prod_j \pi_{1dj}^{x_{dj}}$$

raise to power x_{dj} , where x_{dj} is 0 or 1

One-hot $x_{dj} = \begin{cases} 1 & \text{if } d^{\text{th}} \text{ feature has } j^{\text{th}} \text{ value} \\ 0 & \text{otherwise} \end{cases}$

Then find parameters $\underline{\pi}_0$ of class 0 and parameters $\underline{\pi}_1$ of class 1 to minimize negated log likelihood

$$\text{arg min}_{\underline{\pi}_0, \underline{\pi}_1} - \sum_n \ln[P(\underline{x}_n | y_n)]$$

Can write $\underline{\pi}_0$ as a "stacked" $(D \times J)$ -dim vector; write \underline{x} as a "stacked" $(D \times J)$ -dim vector.

MLE fit $\underline{\pi}_0 = \sum_{n: y_n=0} \underline{x}_n / N_0$ (# examples with $y_n=0$)

$$\underline{\pi}_1 = \sum_{n: y_n=1} \underline{x}_n / N_1$$

Discriminative

train very quickly

- ⊗ interpretable $\exp[\underline{w}_k^T \underline{x}]$
- ⊗ good as "softmax" $\frac{\exp[\underline{w}_k^T \underline{x}]}{\sum_l \exp[\underline{w}_l^T \underline{x}]}$ data sets

Decision boundary From (1) this is \underline{x} s.t.

$$\ln[p(\underline{x}|y=1)] - \ln[p(\underline{x}|y=0)] = \text{constant}$$

For categorical model, this is

$$\ln \left[\prod_d \prod_j \pi_{1dj}^{x_{dj}} \right] - \ln \left[\prod_d \prod_j \pi_{0dj}^{x_{dj}} \right] = \text{constant}$$

$$\Leftrightarrow \sum_d \sum_j x_{dj} \ln \left(\frac{\pi_{1dj}}{\pi_{0dj}} \right) = \text{constant}$$

$$\Leftrightarrow \underline{x}^T \ln \left[\frac{\underline{\pi}_1}{\underline{\pi}_2} \right] = \text{constant}$$

(writing these as stacked $D \times J$ -dim vectors)

Conclude that decision boundary of categorical classifier model is linear. \mathbb{R}

Note!

Multiclass generalizations

$$C = \{C_1, \dots, C_K\}$$

Generatively \rightarrow easy, just use

a categorical class prior (not Bernoulli)
and estimate class conditionals in same way

Classify as $\arg \max_k p(\underline{x} | y = C_k) p(y = C_k)$

Discriminative

Now have ^{separate} parameters \underline{w}_k for each class

$$p(y = C_k | \underline{x}) = \frac{\exp\left[\frac{\underline{w}_k^T \underline{x}}{z}\right]}{\sum_{l=1}^K \exp\left[\frac{\underline{w}_l^T \underline{x}}{z}\right]}$$

"softmax" sums to 1

"Multiclass logistic regression"