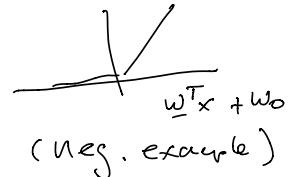


Recall

$$\hat{y}_n = \begin{cases} +1 & \text{if } \underline{w}^\top \underline{x} + w_0 \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

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

of 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

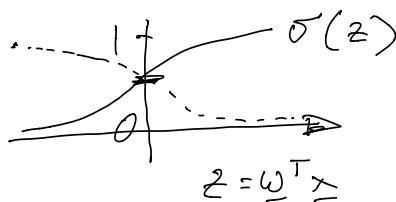
① Convenient to adopt $\{0, 1\}$

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)

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



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

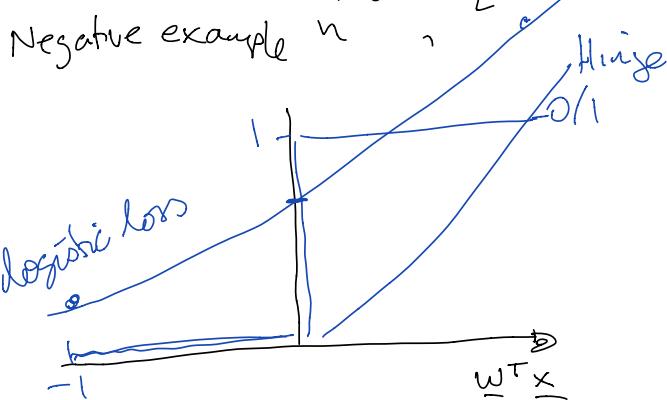
$$\text{Note } p(y=1|x) = \frac{1}{1+e^{-h}} \quad p(y=0|x) = \frac{1}{1+e^h}$$

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

$$\begin{aligned} L_1(\underline{\omega}) &= - \sum_n \ln[p(y_n|x_n)] \\ &= - \sum_n y_n \ln[\sigma(h_n)] - \sum_n (1-y_n) \ln[\sigma(h_n)] \\ &= \underbrace{\sum_n y_n \ln[1 + \exp(-h_n)]}_{\text{Hinge loss}} + \underbrace{\sum_n (1-y_n) \ln[1 + \exp(h_n)]}_{\text{loss on a -ve ex}} \end{aligned}$$



Differentiable \Rightarrow convex

Prefers to make better decisions or correct predictions

Stochastic gradient descent

Positive example

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

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

Negative example

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

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

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

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

Algorithms come from loss function

Loss + SGD \rightarrow Algorithm

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

SGD

[2] Generative approach

Model $p(x, y)$ directly

(1) Decide on a generating process

(2) Parametrize model

(3) Minimize ~~to~~ negative log. likelihood

Use the chain (product) rule:

$$p(x, y) = p(x|y)p(y) \quad y \rightarrow x$$

↓ ↓ ↓
 class conditional class prior
 continuous discrete Bernoulli
 Gaussian Naive Bayes

$$\begin{aligned}
 \boxed{\text{MLE}} \quad & \arg \max_{\underline{w}} \prod_n p(x_n | y_n) = \arg \max_{\underline{w}} \prod_n \frac{p(x_n | y_n)}{p(y_n)} \\
 &= \arg \min_{\underline{w}} - \sum_n \underbrace{\ln [p(x_n | y_n)]}_{\text{solve for parts}} - \sum_n \underbrace{\ln [p(y_n)]}_{\text{solve for parts}}
 \end{aligned}$$

of \underline{w}^* corresponding to class conditional of \underline{w}^* corresponding to class prior

Aside

Classify via Bayes rule

$$p(y=1 | \mathbf{x}) \propto p(y=1) p(\mathbf{x}|y=1)$$

estimate
Class prior

Model as Bernoulli r.v. with

prob θ

$$p(y) = \theta^y (1-\theta)^{1-y} \Rightarrow \frac{\partial \ell(\theta)}{\partial \theta} = 0$$

~~MLE~~ MLE : $\hat{\theta}^* = \sum_n y_n / N$

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

:

solve analytically

Class conditional (continuous)

$p(x|y) \propto$ continuous

$$x|y=0 \sim N(\mu_0, \Sigma_0) \quad x|y=1 \sim N(\mu_1, \Sigma_1)$$

$$\underline{\omega} = \{\mu_0, \Sigma_0, \mu_1, \Sigma_1\} (+ \theta)$$

Intuition: MLE separates by class

For class $y=0$, use

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

\hookrightarrow estimate μ_0, Σ_0

For class $y=1$, use

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

\hookrightarrow estimate μ_1, Σ_1

Example

$$\text{MLE } \mu_0 = \frac{\sum x_n}{N_0} \quad \begin{matrix} \text{---} \\ \text{# examples} \end{matrix} \quad \begin{matrix} \text{---} \\ \text{with } y=0 \end{matrix}$$

\hookrightarrow Understand decision boundaries!

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

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

Boundary:

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

(1) Equivalently, decision boundary $\approx S$:

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

Plug in

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

& 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

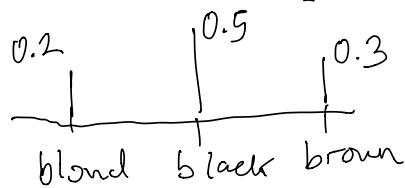
Key assumption (Naive Bayes)

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"

$$\pi_{k d j} \geq 0$$

Probability in class k

of feature d taking on
value j

$$\sum_j \pi_{k d j} = 1$$

for all k ,
all d

Write x_d as a "one-hot" vector

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

Discriminative
Separate parameters w_k & b_k

$$p(y=c_k | \underline{x}) = \frac{\exp[\underline{w}^T \underline{x}]}{\sum_{j=1}^K \exp[\underline{w}_j^T \underline{x}]} = \prod_{d=1}^D \prod_{j=1}^{c_k} \prod_{j \neq c_k}^{c_k} x_{dj}^{x_{dj}}$$

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

Oncer that sum to 1

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

Then find parameters $\underline{\pi}_0$ of class

0 and parameters $\underline{\pi}_1$ of class 1

to minimize negated log likelihood

$$\arg \min_{\underline{\pi}_0, \underline{\pi}_1} - \sum_n \ln [p(\underline{x}_n | y_n)]$$

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

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

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

✓ train very quickly

✓ interpretable

✓ good on small data sets

Decision boundary From (1) this is s.t.

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

For categorical model, this is

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

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

$$\Leftrightarrow \underline{x}^T \ln \begin{bmatrix} \underline{\pi}_1 \\ \underline{\pi}_2 \end{bmatrix} = \text{constant}$$

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

Conclude that decision boundary of categorical classifier model is linear

Note!

Multiclass generalizations

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

Generatively \rightarrow easy, just use

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

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

Discriminative

Now have ^{separate} parameters w_k for each class

$$p(y=C_k | x) = \frac{\exp \left[w_k^T x \right]}{\sum_{l=1}^K \exp \left[w_l^T x \right]}$$

"softmax" sums to 1

"Multiclass Logistic regression"