# PROBABILISTIC GENERATIVE MODELS

$p(x \mid C_k)$   class-conditional densities

$p(C_k)$   class priors

posterior probability   $p(C_k \mid x)$

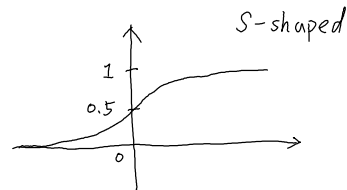for two classes

$$p(C_1 \mid x) = \frac{p(x \mid C_1)\, p(C_1)}{p(x \mid C_1)\, p(C_1) + p(x \mid C_2)\, p(C_2)}$$

$$= \frac{1}{1 + \exp(-a)} = \sigma(a)$$

$$a = \ln \frac{p(x \mid C_1)\, p(C_1)}{p(x \mid C_2)\, p(C_2)}$$

$$\boxed{\sigma(a) = \frac{1}{1 + \exp(-a)} \qquad \text{logistic sigmoid function}}$$

$$\sigma(-a) = 1 - \sigma(a)$$

S-shaped



$$a = \ln \left( \frac{\sigma}{1 - \sigma} \right) \qquad \text{logit function}$$

---

for $K > 2$

$$p(C_k \mid x) = \frac{p(x \mid C_k)\, p(C_k)}{\sum_j p(x \mid C_j)\, p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \ln\, p(x \mid C_k)\, p(C_k)$$

$\boxed{\text{softmax function}}$

if $a_k \gg a_j$ for all $j \neq k$

then   $p(C_k \mid x) \simeq 1$, and   $p(C_j \mid x) \simeq 0$

---

Continuous Inputs   ( Gaussian )

class-conditional densities

$$p(x \mid C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \right\}$$

↑ shared covariance

two-class case :

$$p(C_1 \mid x) = \sigma(w^T x + w_0)$$

$$W = \Sigma^{-1} (\mu_1 - \mu_2)$$

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)}$$
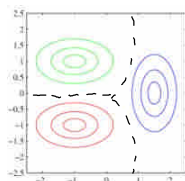


$w_0$

for K classes

$$a_k(x) = W_k^T x + W_{k0}$$

$$W_k = \Sigma^{-1} \mu_k$$

$$w_{k0} = -\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \ln p(C_k)$$

if $p(x|C_k)$ has its own covariance matrix $\Sigma_k$

$\Rightarrow$ quadratic discriminant



How to choose $\mu_k$ for $p(x|C_k)$ ?

$\hookrightarrow$ maximum likelihood solution

parameters:
$\pi, \mu_1, \mu_2, \Sigma$

$t_n = 1$ for $C_1$ , $t_n = 0$ for $C_2$

prior class probality $p(C_1) = \pi$      $p(C_2) = 1 - \pi$

$$p(x_n, C_1) = p(C_1) p(x_n|C_1) = \pi \mathcal{N}(x_n|\mu_1, \Sigma)$$

$$p(x_n, C_2) = p(C_2) p(x_n|C_2) = (1-\pi)\mathcal{N}(x_n|\mu_2, \Sigma)$$

$$p(t | \pi, \mu_1, \mu_2, \Sigma)$$

$$= \prod_{n=1}^{N} \left[ \pi \mathcal{N}(x_n|\mu_1, \Sigma) \right]^{t_n} \left[ (1-\pi) \mathcal{N}(x_n|\mu_2, \Sigma) \right]^{1-t_n}$$

---

maximization w.r.t. $\pi$

terms in the log likelihood function that depend on $\pi$

$$\sum_{n=1}^{N} \left\{ t_n \ln \pi + (1-t_n) \ln (1-\pi) \right\}$$

set the derivative to 0 $\Rightarrow$ $\pi = \frac{1}{N} \sum_{n=1}^{N} t_n = \frac{N_1}{N} = \frac{N_1}{N_1+N_2}$

maximization w.r.t. $\mu_1$

$$\sum_{n=1}^{N} t_n \ln \mathcal{N}(x_n|\mu_1, \Sigma) = -\frac{1}{2} \sum_{n=1}^{N} t_n (x-\mu_1)^T \Sigma^{-1}(x_n-\mu_1) + const.$$

set the derivative to 0 $\Rightarrow$ $\mu_1 = \frac{1}{N_1} \sum_{n=1}^{N} t_n x_n$

similarly $\mu_2 = \frac{1}{N_2} \sum_{n=1}^{N} (1-t_n) x_n$

maximization w.r.t. $\Sigma$

related terms $-\frac{1}{2} \sum_{n=1}^{N} t_n \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^{N} t_n (x_n-\mu_1)^T \Sigma^{-1}(x_n-\mu_1)$

$$-\frac{1}{2} \sum_{n=1}^{N} (1-t_n) \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (1-t_n)(x_n-\mu_2)^T \Sigma^{-1}(x_n-\mu_2)$$

$$= -\frac{N}{2} \ln|\Sigma| - \frac{N}{2} T_r \{\Sigma^{-1}S\}$$

where
$$S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$
$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n-\mu_1)(x_n-\mu_1)^T$$
$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x_n-\mu_2)(x_n-\mu_2)^T$$

Set the derivative with respect to $\Sigma$ to zero
we obtain $\Sigma = S$

$$\left.\begin{array}{l} \frac{\partial}{\partial \Sigma} \ln |\Sigma| = \left(\Sigma^{-1}\right)^T \\[2mm] \frac{\partial}{\partial \Sigma} \mathrm{Tr}\{\Sigma^{-1} S\} = (-\Sigma^{-2})^T S \end{array}\right\} \Rightarrow \quad \left(\Sigma^{-1}\right)^T = \left(\Sigma^{-2}\right)^T S$$

---

Discrete Features

$x_i \in \{0, 1\}$      Naive Bayes
                feature values are treated as
                independent conditioned on the class $C_k$

$$p(\mathbf{x} \mid C_k) = \prod_{i=1}^{D} \mu_{ki}^{x_i} (1-\mu_{ki})^{1-x_i} \quad \rightarrow \quad a_k = \ln p(\mathbf{x}\mid C_k) p(C_k)$$

$$a_k(\mathbf{x}) = \sum_{i=1}^{D} \left\{ x_i \ln \mu_{ki} + (1-x_i) \ln (1-\mu_{ki}) \right\} + \ln p(C_k)$$

linear functions of the input values $x_i$

---

# PROBABILISTIC DISCRIMINATIVE MODELS

difference between probabilistic generative models?
  $\Rightarrow$ indirect approach to finding the parameters of
a generalized linear model, by fitting
class-conditional densities and class priors separately
and then applying Bayes' theorem.

Probabilistic discriminative models:
  direct approach: maximizing a likelihood function
  defined through the conditional distribution $p(C_k \mid \mathbf{x})$.
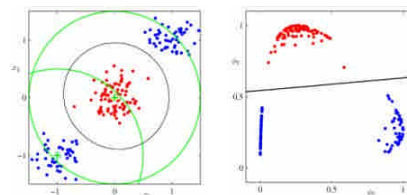
An advantage of the discriminative approach:
fewer parameters to be determined
$\Rightarrow$ improved predictive performance

---

Using fixed basis function again



$\mathbf{x} \rightarrow \phi(\mathbf{x})$

$\phi_0(\mathbf{x}) = 1$

---

LOGISTIC REGRESSION    (for classification not regression)

the posterior probability of class $C_1$ can be written
as a logistic sigmoid acting on a linear function
of the feature vector $\phi$ :

$$p(C_1 \mid \phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

$$p(C_2 | \phi) = 1 - p(C_1 | \phi)$$

$\sigma(\cdot)$ is the logistic sigmoid function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

fewer parameters
$\begin{cases} M\text{-dim } \phi \longrightarrow M \text{ parameters} \\ \text{For generative models, we need } 2M \text{ parameters} \\ \text{for means and } M(M+1)/2 \text{ parameters for} \\ \text{covariance matrix} \end{cases}$

Use maximum likelihood to determine the parameters of the logistic regression model.

To begin with, we write $\quad \dfrac{d\sigma}{da} = \sigma(1-\sigma)$

$\begin{cases} \sigma(a) = \dfrac{1}{1+e^{-a}} \\ \dfrac{\partial \sigma}{\partial a} = \dfrac{1}{(1+e^{-a})^2} \cdot (e^{-a}) = \dfrac{1}{1+e^{-a}} \cdot \dfrac{e^{-a}}{1+e^{-a}} = \sigma(1-\sigma) \end{cases}$

for a data set $\{\phi_n, t_n\}$, where $t_n \in \{0,1\}$, $\phi_n = \phi(x_n)$
with $n = 1, \ldots, N$
the likelihood can be written as

$$p(t | w) = \prod_{n=1}^{N} y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

$$y_n = p(C_1 | \phi_n)$$

define an error function by taking the negative logarithm of the likelihood

$$E(w) = -\ln p(t | w)$$

$$= -\sum_{n=1}^{N} \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}$$

where $\quad y_n = \sigma(a_n), \quad a_n = w^T \phi_n$.

$$\nabla E(w) = \sum_{n=1}^{N} (y_n - t_n) \phi_n$$

$\begin{cases} \text{using} \quad \dfrac{\partial \sigma}{\partial a} = (1-\sigma)\sigma \\ \nabla E(w) = -\sum\limits_{n=1}^{N} \left\{ \dfrac{t_n}{y_n}(y_n(1-y_n)) \phi_n - \dfrac{1-t_n}{1-y_n}(1-y_n)y_n \phi_n \right\} \\ \quad = -\sum\limits_{n=1}^{N} \{t_n(1-y_n)\phi_n - (1-t_n)y_n \phi_n\} \end{cases}$

So the contribution to the gradient from data point $n$ is given by the error $y_n - t_n$ times the basis function vector $\phi_n$

$\left(\begin{array}{c} \text{no closed-form} \\ \text{solution} \end{array}\right)$

We may derive a sequential-update algorithm

$$w^{(\tau+1)} = w^{(\tau)} - \eta (y_n - t_n) \phi_n$$
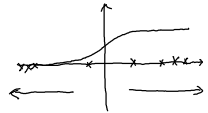
maximum likelihood
→ overfitting for data sets that are linearly separable

multiple solutions    which one is the best ?
we need to include a prior for MAP. or add
a regularization term

$w^T \phi = 0$    $\sigma = 0.5$
decision plane

$w^T \phi_n \geqslant 0$  $t_n = 1$  ,    $w^T \phi_n < 0$  $t_n = 0$

⇓

Heaviside step function

$$\boxed{\text{ITERATIVE REWEIGHTED LEAST SQUARES}}$$

Recall:  linear regression , maximum likelihood solution
   assume a Gaussian noise model
   → closed-form solution

For logistic regression , no closed-form solution

$$\nabla E(w) = \sum_{n=1}^{N} (y_n - t_n) \phi_n \overset{?}{=} 0$$

$$y_n = \sigma (w^T \phi_n)  \qquad \underset{n}{\Sigma} \sigma$$
$$\sigma^{-1} \quad ?$$

---

NEWTON - RAPHSON ITERATIVE OPTIMIZATION

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w^{(old)})$$

H is the Hessian matrix whose elements comprise
the second derivatives of $E(w)$ w.r.t. w

$$\boxed{H_{ij} = \frac{\partial^2 E}{\partial w_i \, \partial w_j}}$$

1D   Newton's method
   $x^{(new)} = x^{(old)} - g(x^{(old)}) / g'(x^{(old)})$    finding $g = 0$
   $x^{(new)} = x^{(old)} - f'(x^{(old)}) / f''(x^{(old)})$
      finding $f' = 0$

Try to apply the Newton-Raphson method to
linear regression as a practice

$$\nabla E(w) = \sum_{n=1}^{N} (w^T \phi_n - t_n) \phi_n = \bar{\Phi}^T \bar{\Phi} w - \bar{\Phi}^T t$$

$$H = \nabla \nabla E(w) = \sum_{n=1}^{N} \phi_n \phi_n^T = \bar{\Phi}^T \bar{\Phi}$$

Newton's update

$$w^{(new)} = w^{(old)} - (\bar{\Phi}^T \bar{\Phi})^{-1} \{ \bar{\Phi}^T \bar{\Phi} w^{(old)} - \bar{\Phi}^T t \}$$

$$= (\bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T t$$

we get the standard least-squares solution

$$\bar{\Phi} = \begin{bmatrix} \vdots \\ \phi_n^T \\ \vdots \end{bmatrix}$$
$$N \times M$$

Now try to apply the Newton-Raphson method to

$$E(w) = -\ln p(t \mid w) = -\sum_{n=1}^{N} \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}$$

$$\nabla E(w) = \sum_{n=1}^{N} (y_n - t_n) \phi_n = \Phi^T (y - t)$$

$$H = \nabla\nabla E(w) = \sum_{n=1}^{N} y_n (1-y_n) \phi_n \phi_n^T = \Phi^T R \Phi$$

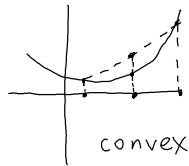The Hessian is no longer constant; it depends on $w$ through the weighting matrix $R$.

$R$ $N \times N$ diagonal
$R_{nn} = y_n(1-y_n)$

$$0 < y_n < 1$$
$$\Rightarrow \quad v^T H v > 0 \quad \text{for an arbitrary } v$$
$$\Rightarrow \quad H \text{ is positive definite}$$

The error function is a convex function of $w$ and hence has a unique minimum.

$$w^{(new)} = w^{(old)} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$
$$= (\Phi^T R \Phi)^{-1} \{ \Phi^T R \Phi w^{(old)} - \Phi^T (y-t) \}$$
$$= (\Phi^T R \Phi)^{-1} \Phi^T R Z$$

$$Z = \Phi w^{(old)} - R^{-1}(y - t) \qquad \underline{\text{meaning ?}}$$

convex

$R$ depends on $w$. we need to update $w$ iteratively

(Iterative Reweighted Least Squares, IRLS)

---

The weighting matrix $R$ can be interpreted as variances

$$E[t] = \sigma(x) = y$$

Bournelli   P.685
$$p(t) = y^t (1-y)^{1-t}$$

$$Var[t] = E[t^2] - E[t]^2 = E[t] - E[t]^2$$
$$= y - y^2 = y(1-y)$$

$t \in \{0, 1\}$
$t^2 = 1$

linearized problem in the space of $a = w^T \phi$

local linear approximation to logistic function

$$a_n(w) \simeq a_n(w^{(old)}) + \frac{d a_n}{d y_n}\Big|_{w^{(old)}} (t_n - y_n)$$
$$= \phi_n^T w^{(old)} - \frac{(y_n - t_n)}{y_n(1-y_n)} = Z_n$$

$n$-th element of $Z$

Meaning ?

$Z_n$: as an effective target value in the space obtained by making a local linear approximation to the logistic sigmoid function around the current operating point $w^{(old)}$

$$\sigma = \frac{1}{1+e^{-a}}$$
$$a = \ln\left(\frac{\sigma}{1-\sigma}\right)$$
$$\frac{da}{d\sigma} = \frac{1}{\sigma(1-\sigma)}$$
$$\sigma = y \qquad \text{see } (4.61)$$
$$(4.88)$$

(compared with least-squares solutions)

approximate $y_n$ by $a_n$

$a = w\phi$
$\sigma(a)$