

Logistic Regression

- Two classes +1, -1
- Model $p(y|x)$ (rather than $p(x|y)$ as in LDA)

$$p(1|x) = \frac{\exp(\beta_0 + \beta^T x)}{1 + \exp(\beta_0 + \beta^T x)}$$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix}$$

$$p(-1|x) = \frac{1}{1 + \exp(\beta_0 + \beta^T x)} = 1 - p(1|x)$$

$g(t) = \frac{e^t}{1+e^t}$ is "logistic" function

Estimating Parameters

Given training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

find β_0 and β that maximize

$$h(\beta_0, \beta) = \prod_{i=1}^N p(y_i|x_i)$$

Maximum likelihood estimation - done using iterative optimization methods

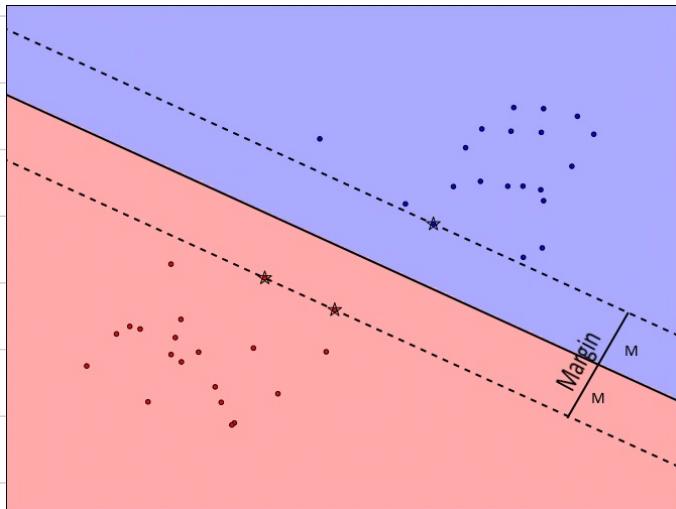
Logistic Regression is a linear classifier:

$$\begin{aligned} p(1|x) &\stackrel{1}{\geq} p(-1|x) \equiv \ln p(1|x) \stackrel{1}{\geq} \ln p(-1|x) \\ &\equiv \hat{\beta}_0 + \hat{\beta}^T x \stackrel{1}{\geq} 0 \end{aligned}$$

linear discriminant function

Support Vector Machine (SVM)

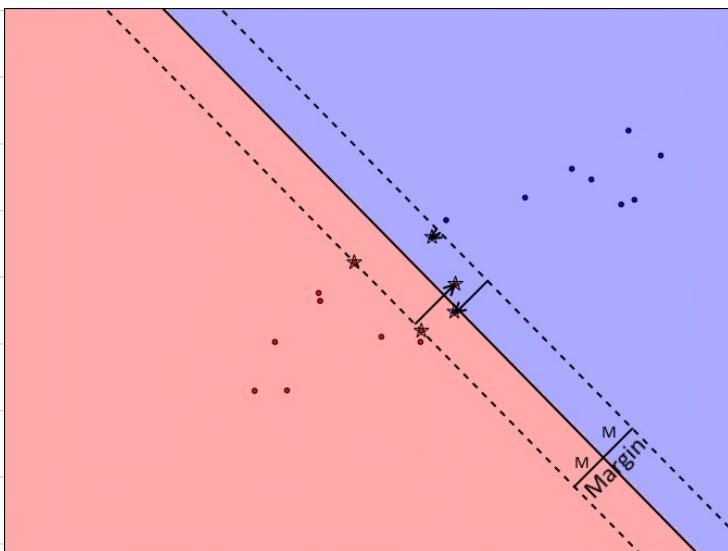
Linearly separable Training Data



Ininitely many
linear boundaries
separating classes

SVM maximizes
margin M between
classes

Training Data Not Linearly Separable



1. Fix decision boundary
and margin m

2. Compute distances of
points on "wrong"
side of margin to
margin boundary ξ_i^*

3. Compute penalty

$$\xi_i = \frac{\xi_i^*}{m}$$

4. Maximize m subject to $\sum_i \xi_i \leq C$ to get M

5. Choose boundary to maximize M

Naive Bayes Classifier

- Recall $\hat{y}_{\text{Bayes}} = \arg \max_{y=1, \dots, M} \pi_y p(x|y)$

Feature vector $\underline{x} = [x_1, x_2, \dots, x_d]^T$

- Naive Bayes assumes that x_1, x_2, \dots, x_d are independent, conditioned on label y , i.e.

$$p(x|y) = p(x_1|y) \cdot p(x_2|y) \cdots p(x_d|y)$$

- Get estimates $\hat{\pi}_y, \hat{p}(x_1|y), \hat{p}(x_2|y) \cdots \hat{p}(x_d|y)$

for each $y = 1, 2, \dots, M$. Then,

$$\begin{aligned}\hat{y}_{\text{NB}} &= \arg \max_{y=1, 2, \dots, M} \hat{\pi}_y \prod_{j=1}^d \hat{p}(x_j|y) \\ &= \arg \max_{y=1, 2, \dots, M} \ln \hat{\pi}_y + \sum_{j=1}^d \ln \hat{p}(x_j|y)\end{aligned}$$

Pros and Cons of NB classifier

+ estimating $p(x|y)$ much harder than estimating $p(x_1|y), p(x_2|y), \dots, p(x_d|y)$ separately

+ can mix different kinds of features:
continuous real-valued, discrete, categorical

- independence assumption is simplistic

Finding $\hat{p}(x_j | y)$ using training data

Convenient to use parametric model:

$$p(x_j | y) = p(x_j | y, \theta_{jy})$$

Then, $\hat{p}(x_j | y) = p(x_j | y, \hat{\theta}_{jy})$

Example (continuous x_j)

$$p(x_j | y, \theta_{jy}) = \frac{1}{\sqrt{2\pi \sigma_{jy}^2}} e^{-\frac{(x_j - \mu_{jy})^2}{2\sigma_{jy}^2}} \leftarrow N(\mu_{jy}, \sigma_{jy}^2)$$

$$\theta_{jy} = (\mu_{jy}, \sigma_{jy}^2)$$

Given training data $\{(x_i, y_i)\}_{i=1}^N$,

$$\hat{\mu}_{jy} = \frac{\sum_{i:y_i=y} x_{ij}}{N_y}, \quad \hat{\sigma}_{jy}^2 = \frac{\sum_{i:y_i=y} (x_{ij} - \hat{\mu}_{jy})^2}{N_y} \leftarrow \text{or } N_y - 1$$

Example (discrete or categorical x_j)

x_j takes values in set $\{a_1, a_2, \dots, a_e\}$ with probabilities $s_{1,jy}, s_{2,jy}, \dots, s_{e,jy}$, for class y .

$$\theta_{jy} = (s_{1,jy}, s_{2,jy}, \dots, s_{e,jy}), \quad \hat{\theta}_{jy} = (\hat{s}_{1,jy}, \hat{s}_{2,jy}, \dots, \hat{s}_{e,jy})$$

$$\hat{s}_{k,jy} = \frac{\# \text{ training data in class } y \text{ with } x_{ij} = a_k}{N_y}$$

(See typed notes for Laplacian smoothing so $\hat{s} \neq 0$)