

6. Probabilistic models

References

C. Bishop. *Pattern Recognition and Machine Learning*. Sect. 4.2, 4.3

Approaches that allows to estimate probability that given an instance we have a certain classes.

Two families of models:

Generative: estimate $P(\mathbf{x}|C_i)$ and then compute $P(C_i|\mathbf{x})$ with Bayes

Discriminative: estimate $P(C_i|\mathbf{x})$ directly

6.1 Probabilistic generative models

$$\begin{aligned} P(C_1|\mathbf{x}) &= \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x}|C_1)P(C_1) + P(\mathbf{x}|C_2)P(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

with:

$$a = \ln \frac{p(\mathbf{x}|C_1)P(C_1)}{p(\mathbf{x}|C_2)P(C_2)}$$

and

$$\sigma(a) = \frac{1}{1+\exp(-a)} \text{ the sigmoid function.}$$

Multi-class

$$P(C_k|\mathbf{x}) = \frac{P(\mathbf{x}|C_k)P(C_k)}{\sum_j P(\mathbf{x}|C_j)P(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

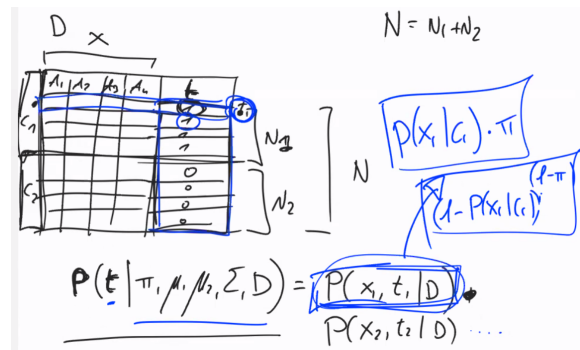
(normalized exponential or softmax function)

with $a_k = \ln P(\mathbf{x}|C_k)P(C_k)$

6.1.1 Maximum likelihood

Likelihood function

$$P(\mathbf{t}|\pi, \mu_1, \mu_2, \Sigma, D) = \prod_{n=1}^N \underbrace{[\pi \mathcal{N}(\mathbf{x}_n; \mu_1, \Sigma)]^{t_n}}_{P(\mathbf{x}_n|C_1)} \underbrace{[(1-\pi) \mathcal{N}(\mathbf{x}_n; \mu_2, \Sigma)]^{1-t_n}}_{P(\mathbf{x}_n|C_2)}$$



For 2 classes, we obtain

$$\pi = \frac{N_1}{N} \rightarrow \text{class 1.}$$

$$\pi = \frac{N_2}{N} \rightarrow \text{Total}$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$\text{with } S_i = \frac{1}{N_i} \sum_{n \in C_i} (\mathbf{x}_n - \mu_i)(\mathbf{x}_n - \mu_i)^T, i = 1, 2$$

Posterior distributions with parametric models:

For two classes

$$P(C_1 | \mathbf{x}) = \sigma(a)$$

For $k \geq 2$ classes

$$P(C_i | \mathbf{x}) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \mathbf{w}^T \mathbf{x} + w_0$$

$$\mathbf{w}^T \mathbf{x} + w_0 = \begin{pmatrix} w_0 & \mathbf{w} \end{pmatrix} \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Likelihood for a parametric model \mathcal{M}_Θ : $P(\mathbf{t}|\Theta, D)$, $D = \langle \mathbf{X}, \mathbf{t} \rangle$

Maximum likelihood solution:

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \ln P(\mathbf{t}|\Theta, \mathbf{X})$$

$$\tilde{\mathbf{w}} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}, \tilde{\mathbf{x}} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

When \mathcal{M}_Θ belongs to the exponential family, likelihood $P(\mathbf{t}|\Theta, \mathbf{X})$ can be expressed in the form $P(\mathbf{t}|\tilde{\mathbf{w}}, \mathbf{X})$, with maximum likelihood

$$\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}}}{\operatorname{argmax}} \ln P(\mathbf{t}|\tilde{\mathbf{w}}, \mathbf{X})$$

$$a_k = \mathbf{w}^T \mathbf{x} + w_0 = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

6.2 Probabilistic discriminative models

Estimate directly

$$P(C_i|\tilde{\mathbf{x}}, D) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

with maximum likelihood

$$\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}}}{\operatorname{argmax}} \ln P(\mathbf{t}|\tilde{\mathbf{w}}, \mathbf{X})$$

max this likelihood

without estimating the model parameters.

Simplified notation (dataset omitted): $P(\mathbf{t}|\tilde{\mathbf{w}})$

6.2.1 Logistic regression

For 2 classes, likelihood function:

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

Note: t_n : value in the data set
corresponding to \mathbf{x}_n , y_n : posterior
prediction of the current model
for \mathbf{x}_n .

$$\text{with } y_n = p(C_1|\mathbf{x}_n) = \sigma(\mathbf{w}^T \mathbf{x}_n)$$

$$E(\mathbf{w}) \equiv -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)]$$

Solution concept: solve the optimization problem

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w})$$

Cross-entropy error function:

To minimize error we apply Newton-Raphson:

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

Gradient descent step

$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

$\mathbf{H} = \nabla \nabla E(\mathbf{w})$ is the Hessian matrix of $E(\mathbf{w})$ (second derivatives with respect to \mathbf{w}).

Given

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \dots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix} \quad \mathbf{t} = \begin{pmatrix} t_1 \\ \dots \\ t_N \end{pmatrix},$$

$\mathbf{y}(\tilde{\mathbf{w}}) = (y_1, \dots, y_N)^T$ posterior predictions of model $\tilde{\mathbf{w}}$

$\mathbf{R}(\tilde{\mathbf{w}})$: diagonal matrix with $R_{nn} = y_n(1 - y_n)$

we have

$$\nabla E(\tilde{\mathbf{w}}) = \tilde{\mathbf{X}}^T (\mathbf{y}(\tilde{\mathbf{w}}) - \mathbf{t})$$

$$\mathbf{H}(\tilde{\mathbf{w}}) = \nabla \nabla E(\tilde{\mathbf{w}}) = \sum_{n=1}^N y_n(1 - y_n) \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T = \tilde{\mathbf{X}}^T \mathbf{R}(\tilde{\mathbf{w}}) \tilde{\mathbf{X}}$$

Iterative method:

1. Initialize \mathbf{w}
2. Repeat until termination condition

$$\mathbf{w} \leftarrow \mathbf{w} - (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{t})$$

Multiclass:

K classes

$$P(C_k|\tilde{\mathbf{x}}) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}} \quad k = 1, \dots, K$$

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \dots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix} \quad \mathbf{T} = \begin{pmatrix} t_1^T \\ \dots \\ t_N^T \end{pmatrix} \quad \text{1-of-}K \text{ encoding of labels}$$

$\mathbf{y}_n^T = (y_{n1} \dots y_{nK})^T$ posterior prediction of $\tilde{\mathbf{x}}_n$ for model $\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K$

$$\mathbf{Y}(\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = \begin{pmatrix} \mathbf{y}_1^T \\ \dots \\ \mathbf{y}_N^T \end{pmatrix} \quad \text{posterior predictions of model } \tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K$$

Discriminative model

$$P(\mathbf{T}|\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = \prod_{n=1}^N \prod_{k=1}^K P(C_k|\tilde{\mathbf{x}}_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

with $y_{nk} = \mathbf{Y}[n, k]$ and $t_{nk} = \mathbf{T}[n, k]$.

Cross-entropy error function

$$E(\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = -\ln P(\mathbf{T}|\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

Iterative algorithm

gradient $\nabla_{\tilde{\mathbf{w}}_j} E(\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = \dots$

Hessian matrix $\nabla_{\tilde{\mathbf{w}}_k} \nabla_{\tilde{\mathbf{w}}_j} E(\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K) = \dots$

SUMMARY

Given a target function $f : X \rightarrow C$, and data set D

assume a parametric model for the posterior probability $P(C_k|\tilde{\mathbf{x}}, \tilde{\mathbf{w}})$
 $\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$ (2 classes) or $\frac{\exp(\tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}})}{\sum_{j=1}^K \exp(\tilde{\mathbf{w}}_j^T \tilde{\mathbf{x}})}$ (k classes)

Define an error function $E(\tilde{\mathbf{w}})$ (negative log likelihood)

Solve the optimization problem

$$\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} E(\tilde{\mathbf{w}})$$

Classify new sample $\tilde{\mathbf{x}}'$ as C_{k^*} where $k^* = \operatorname{argmax}_{k=1,\dots,K} P(C_k|\tilde{\mathbf{x}}', \tilde{\mathbf{w}}^*)$