

Discriminative models

Machine Learning

Hamid R Rabiee – Zahra Dehghanian
Spring 2025



Sharif University
of Technology

Probabilistic classifiers

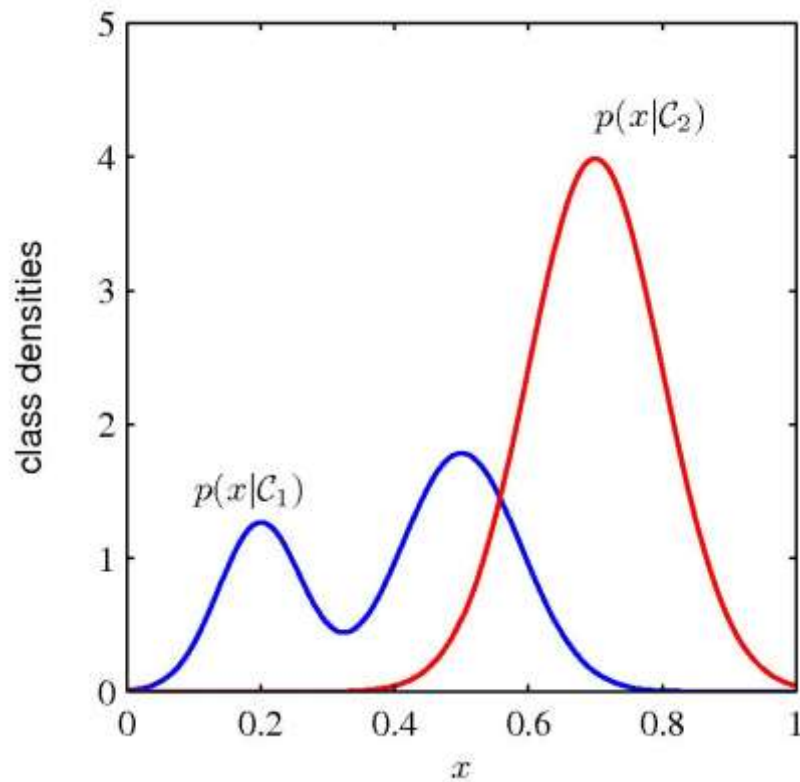
- How can we find the probabilities required in the Bayes decision rule?
- Probabilistic classification approaches can be divided in two main categories:
 - **Generative**
 - Estimate pdf $p(\mathbf{x}, \mathcal{C}_k)$ for each class \mathcal{C}_k and then use it to find $p(\mathcal{C}_k|\mathbf{x})$
 - or alternatively estimate both pdf $p(\mathbf{x}|\mathcal{C}_k)$ and $p(\mathcal{C}_k)$ to find $p(\mathcal{C}_k|\mathbf{x})$
 - **Discriminative**
 - Directly estimate $p(\mathcal{C}_k|\mathbf{x})$ for each class \mathcal{C}_k



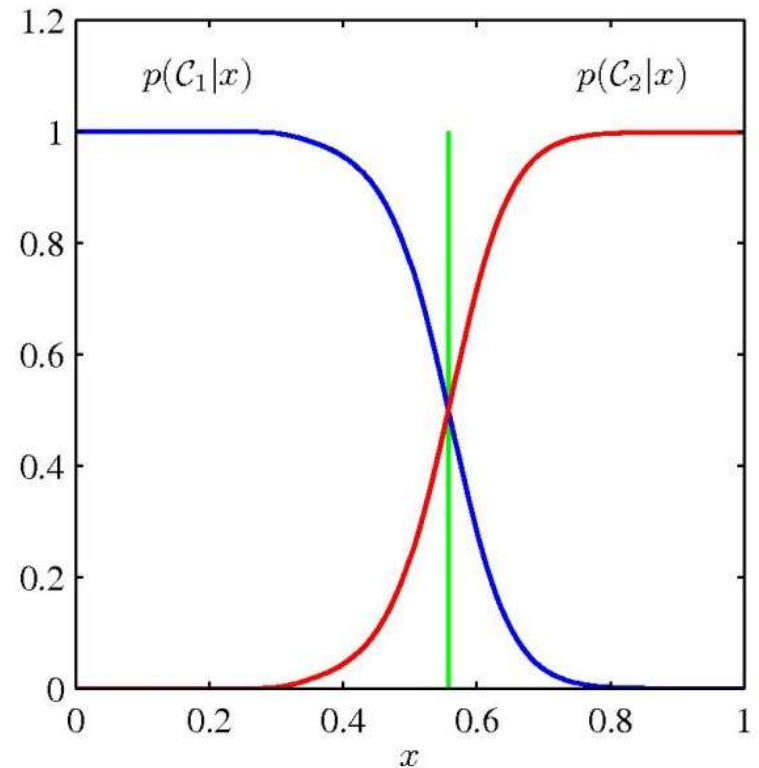
Generative approach

- Inference stage
 - Determine class conditional densities $p(\mathbf{x}|\mathcal{C}_k)$ and priors $p(\mathcal{C}_k)$
 - Use the Bayes theorem to find $p(\mathcal{C}_k|\mathbf{x})$
- Decision stage: After learning the model (inference stage), make optimal class assignment for new input
 - if $p(\mathcal{C}_i|\mathbf{x}) > p(\mathcal{C}_j|\mathbf{x}) \quad \forall j \neq i$ then decide \mathcal{C}_i

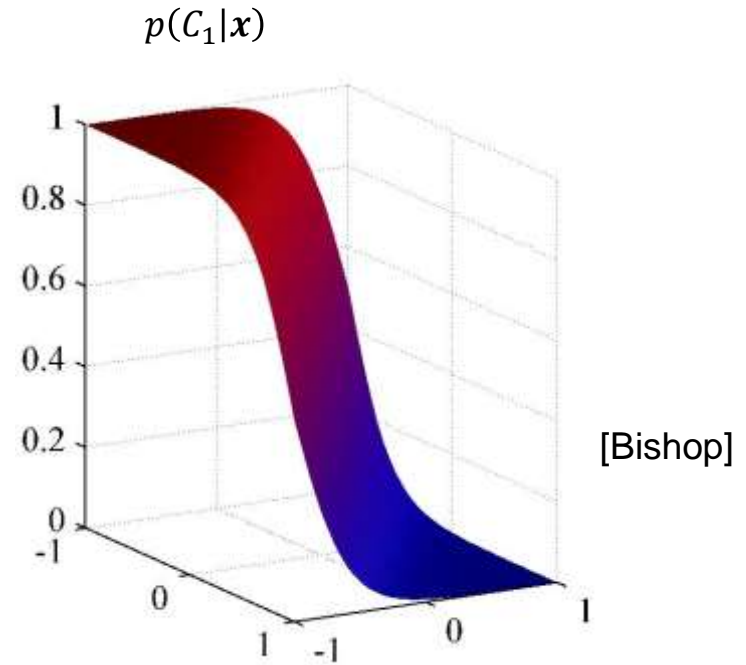
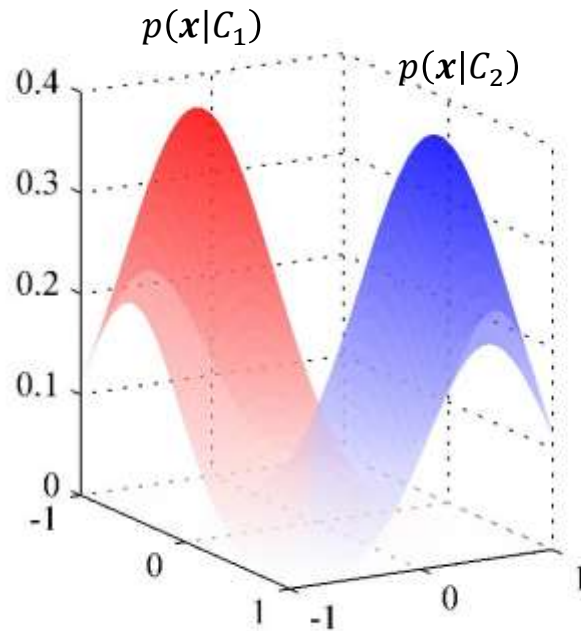
Discriminative vs. generative approach



[Bishop]



Class conditional densities vs. posterior



$$p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$
$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Probabilistic discriminant functions

- **Discriminant functions:** A popular way of representing a classifier
 - A discriminant function $f_i(\mathbf{x})$ for each class \mathcal{C}_i ($i = 1, \dots, K$):
 - \mathbf{x} is assigned to class \mathcal{C}_i if:

$$f_i(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq i$$

- Representing Bayesian classifier using discriminant functions:
 - Classifier minimizing error rate: $f_i(\mathbf{x}) = P(\mathcal{C}_i|\mathbf{x})$
 - Classifier minimizing risk: $f_i(\mathbf{x}) = -\sum_{j=1}^K L_{ij}p(\mathcal{C}_j|\mathbf{x})$

Discriminative approach

- Inference stage
 - Determine the posterior class probabilities $P(C_k|\mathbf{x})$ directly
- Decision stage: After learning the model (inference stage), make optimal class assignment for new input
 - if $P(C_i|\mathbf{x}) > P(C_j|\mathbf{x}) \quad \forall j \neq i$ then decide C_i

Discriminative approach: logistic regression

$K = 2$

- More general than discriminant functions:
 - $f(\mathbf{x}; \mathbf{w})$ predicts posterior probabilities $P(y = 1|\mathbf{x})$

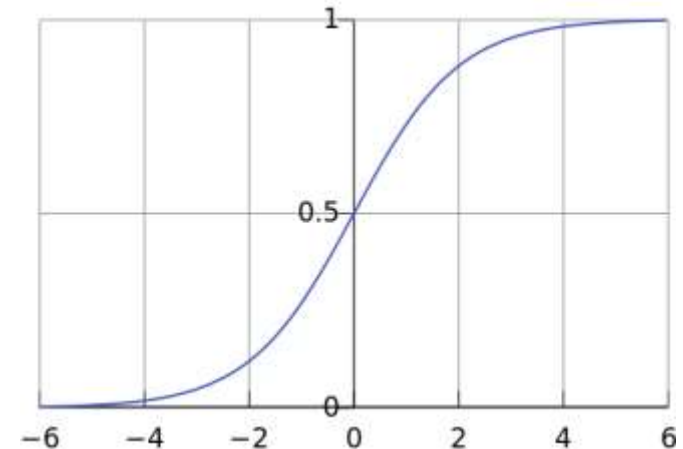
$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$\sigma(\cdot)$ is an activation function

$$\mathbf{x} = [1, x_1, \dots, x_d]$$
$$\mathbf{w} = [w_0, w_1, \dots, w_d]$$

- Sigmoid (logistic) function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Logistic regression

- $f(\mathbf{x}; \mathbf{w})$: probability that $y = 1$ given \mathbf{x} (parameterized by \mathbf{w})

$$P(y = 1|\mathbf{x}; \mathbf{w}) = f(\mathbf{x}; \mathbf{w})$$

$K = 2$
 $y \in \{0,1\}$

$$P(y = 0|\mathbf{x}; \mathbf{w}) = 1 - f(\mathbf{x}; \mathbf{w})$$

$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$0 \leq f(\mathbf{x}; \mathbf{w}) \leq 1$$

estimated probability of $y = 1$ on input \mathbf{x}

- ▶ Example: Cancer (Malignant, Benign)

- ▶ $f(\mathbf{x}; \mathbf{w}) = 0.7$

- ▶ 70% chance of tumor being malignant

Logistic regression: Decision surface

- Decision surface $f(\mathbf{x}; \mathbf{w}) = \text{constant}$

- $f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x})}} = 0.5$

- Decision surfaces are linear functions of \mathbf{x}

if $f(\mathbf{x}; \mathbf{w}) \geq 0.5$ then $y = 1$

else $y = 0$

Equivalent to

if $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$ then $y = 1$

else $y = 0$

Logistic regression: ML estimation

- Maximum (conditional) log likelihood:

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \log \prod_{i=1}^n p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)})$$

$$p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) = f(\mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}} (1 - f(\mathbf{x}^{(i)}; \mathbf{w}))^{(1-y^{(i)})}$$

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \sum_{i=1}^n \left[y^{(i)} \log (f(\mathbf{x}^{(i)}; \mathbf{w})) + (1 - y^{(i)}) \log (1 - f(\mathbf{x}^{(i)}; \mathbf{w})) \right]$$

Logistic regression: cost function

- $$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

$$\begin{aligned} J(\mathbf{w}) &= - \sum_{i=1}^n \log p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n -y^{(i)} \log(f(\mathbf{x}^{(i)}; \mathbf{w})) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}; \mathbf{w})) \end{aligned}$$

- No closed form solution for

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$$

- However $J(\mathbf{w})$ is convex.

Logistic regression: Gradient descent

-

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^t)$$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^n (f(\mathbf{x}^{(i)}; \mathbf{w}) - y^{(i)}) \mathbf{x}^{(i)}$$

► Is it similar to gradient of SSE for linear regression?

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$$

Logistic regression: loss function

-

$$\text{Loss}(y, f(\mathbf{x}; \mathbf{w})) = -y \times \log(f(\mathbf{x}; \mathbf{w})) - (1 - y) \times \log(1 - f(\mathbf{x}; \mathbf{w}))$$

$$\text{Since } y = 1 \text{ or } y = 0 \Rightarrow \text{Loss}(y, f(\mathbf{x}; \mathbf{w})) = \begin{cases} -\log(f(\mathbf{x}; \mathbf{w})) & \text{if } y = 1 \\ -\log(1 - f(\mathbf{x}; \mathbf{w})) & \text{if } y = 0 \end{cases}$$

How is it related to zero-one loss?

$$\text{Loss}(y, \hat{y}) = \begin{cases} 1 & y \neq \hat{y} \\ 0 & y = \hat{y} \end{cases}$$

$$f(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Logistic regression: cost function (summary)

- Logistic Regression (LR) has a more proper cost function for classification than SSE and Perceptron

- Why is the cost function of LR also more suitable than?

$$J(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$

where $f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$

- The conditional distribution $p(y|\mathbf{x}; \mathbf{w})$ in the classification problem is not Gaussian (it is Bernoulli)
- The cost function of LR is convex

Posterior probabilities

- Two-class: $p(\mathcal{C}_k|\mathbf{x})$ can be written as a logistic sigmoid for a wide choice of $p(\mathbf{x}|\mathcal{C}_k)$ distributions

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(a(\mathbf{x})) = \frac{1}{1 + \exp(-a(\mathbf{x}))}$$

- Multi-class: $p(\mathcal{C}_k|\mathbf{x})$ can be written as a soft-max for a wide choice of $p(\mathbf{x}|\mathcal{C}_k)$

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{\exp(a_k(\mathbf{x}))}{\sum_{j=1}^K \exp(a_j(\mathbf{x}))}$$

Multi-class logistic regression

- For each class k , $f_k(\mathbf{x}; \mathbf{W})$ predicts the probability of $y = k$
 - ▶ i.e., $P(y = k|\mathbf{x}, \mathbf{W})$
- ▶ On a new input \mathbf{x} , to make a prediction, pick the class that maximizes $f_k(\mathbf{x}; \mathbf{W})$:

$$\alpha(\mathbf{x}) = \operatorname{argmax}_{k=1,\dots,K} f_k(\mathbf{x})$$

if $f_k(\mathbf{x}) > f_j(\mathbf{x}) \quad \forall j \neq k$ then
decide C_k

Multi-class logistic regression

$$K > 2$$

$$y \in \{1, 2, \dots, K\}$$

$$f_k(\mathbf{x}; \mathbf{W}) = p(y = k | \mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x})}$$

- Normalized exponential (aka softmax)
 - If $\mathbf{w}_k^T \mathbf{x} \gg \mathbf{w}_j^T \mathbf{x}$ for all $j \neq k$ then $p(C_k | \mathbf{x}) \simeq 1, p(C_j | \mathbf{x}) \simeq 0$

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{\sum_{j=1}^K p(\mathbf{x} | C_j) p(C_j)}$$

Logistic regression: multi-class

$$\hat{\mathbf{W}} = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$
$$J(\mathbf{W}) = -\log \prod_{i=1}^n p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \mathbf{W})$$

$$= -\log \prod_{i=1}^n \prod_{k=1}^K f_k(\mathbf{x}^{(i)}; \mathbf{W})^{y_k^{(i)}}$$

$$= -\sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log(f_k(\mathbf{x}^{(i)}; \mathbf{W}))$$

$$\mathbf{W} = [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_K]$$

\mathbf{y} is a vector of length K (1-of- K coding)
e.g., $\mathbf{y} = [0, 0, 1, 0]^T$ when the target class is C_3

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \vdots \\ \mathbf{y}^{(n)} \end{bmatrix} = \begin{bmatrix} y_1^{(1)} & \cdots & y_K^{(1)} \\ \vdots & \ddots & \vdots \\ y_1^{(n)} & \cdots & y_K^{(n)} \end{bmatrix}$$

Logistic regression: multi-class

-

$$\mathbf{w}_j^{t+1} = \mathbf{w}_j^t - \eta \nabla_{\mathbf{W}} J(\mathbf{W}^t)$$

$$\nabla_{\mathbf{w}_j} J(\mathbf{W}) = \sum_{i=1}^n \left(f_j(\mathbf{x}^{(i)}; \mathbf{W}) - y_j^{(i)} \right) \mathbf{x}^{(i)}$$

- We usually consider also a regularization term and the gradient will be

$$\nabla_{\mathbf{w}_j} J(\mathbf{W}) = \lambda \mathbf{W} + \sum_{i=1}^n \left(f_j(\mathbf{x}^{(i)}; \mathbf{W}) - y_j^{(i)} \right) \mathbf{x}^{(i)}$$

Log-odds Ratio

- Optimal rule $y = \arg \max_c p(t = c|x)$ is equivalent to

$$\begin{aligned} y = c &\Leftrightarrow \frac{p(t = c|x)}{p(t = j|x)} \geq 1 \quad \forall j \neq c \\ &\Leftrightarrow \log \frac{p(t = c|x)}{p(t = j|x)} \geq 0 \quad \forall j \neq c \end{aligned}$$

- For the binary case

$$y = 1 \Leftrightarrow \log \frac{p(t = 1|x)}{p(t = 0|x)} \geq 0$$

Logistic Regression (LR): summary

- LR is a linear classifier
- LR optimization problem is obtained by maximum likelihood
 - when assuming Bernoulli distribution for conditional probabilities whose mean is $\frac{1}{1+e^{-(w^T x)}}$
- No closed-form solution for its optimization problem
 - But convex cost function and global optimum can be found by gradient ascent

Discriminative vs. generative: number of parameters

- d -dimensional feature space
- Logistic regression: $d + 1$ parameters
 - $\mathbf{w} = (w_0, w_1, \dots, w_d)$
- Generative approach:
 - Gaussian class-conditionals with shared covariance matrix
 - $2d$ parameters for means
 - $d(d + 1)/2$ parameters for shared covariance matrix
 - one parameter for class prior $p(C_1)$.
- But LR is more robust, less sensitive to incorrect modeling assumptions

Summary of alternatives

- Generative

- ▶ Most demanding, because it finds the joint distribution $p(\mathbf{x}, \mathcal{C}_k)$
- ▶ Usually needs a large training set to find $p(\mathbf{x}|\mathcal{C}_k)$
- ▶ Can find $p(\mathbf{x}) \Rightarrow$ Outlier or novelty detection

- ▶ Discriminative

- ▶ Specifies what is really needed (i.e., $p(\mathcal{C}_k|\mathbf{x})$)
- ▶ More computationally efficient

Feed back

? <https://forms.gle/vKRbyVVsWRKcZuqr8>



Resources

- C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 4.2-4.3.
- Course CE-717, Dr. M.Soleymani

