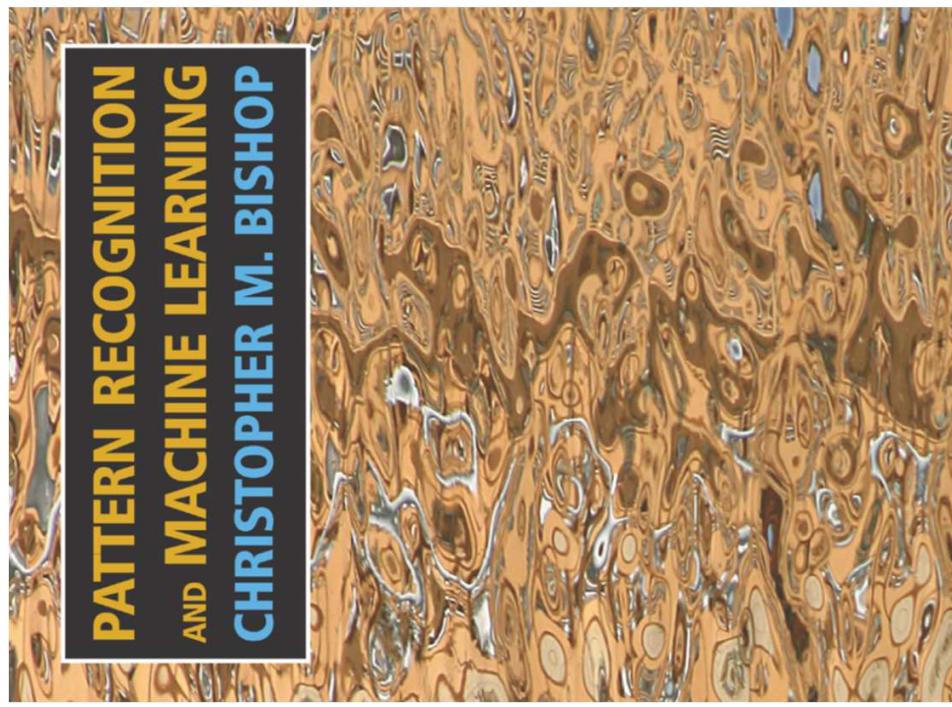


# Discriminative Classifiers

Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk))

# Reference



Today's class roughly follows  
Chapter 4.3

Pattern Recognition and  
Machine Learning

Christopher Bishop, 2006

# Discriminative Classifier

---

- Target: infer  $p(y|\mathbf{x})$  given dataset  $D$ .
  - Step 1. Making a model assumption  $p(y|\mathbf{x}; \mathbf{w})$ .
  - Step 2. Construct the likelihood function  $p(D|\mathbf{w})$ .
  - Step 3. Estimate the parameters: MLE, MAP, Full Prob...
- 
- First Question: What model should we use?
  - MVN? NO, that is for continuous variable.
  - Our output  $y$  is clearly a discrete value.

# Modelling $p(y|\boldsymbol{x})$

---

- Can we express  $p(y|\boldsymbol{x})$  using  $p(\boldsymbol{x}|y)$ ?

- Bayes rule says:

$$\bullet p(y|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|y)p(y)}{p(\boldsymbol{x})} = \frac{p(\boldsymbol{x}|y)p(y)}{\sum_{y'} p(\boldsymbol{x}|y')p(y')} = \frac{p(\boldsymbol{x}|y)p(y)}{\sum_{y'} p(\boldsymbol{x}|y')p(y')}$$

Marginalization!

- Suppose  $y \in \{-1, 1\}$

$$\bullet p(y=1|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|y=1)p(y=1)}{p(\boldsymbol{x}|y'=1)p(y'=1) + p(\boldsymbol{x}|y'=-1)p(y'=-1)}$$

# Modelling $p(y|\boldsymbol{x})$

---

- Suppose  $y \in \{-1, 1\}$

$$\bullet p(y = 1|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|y = 1)p(y = 1)}{p(\boldsymbol{x}|y' = 1)p(y' = 1) + p(\boldsymbol{x}|y' = -1)p(y' = -1)}$$

- Nothing has changed, but we are representing  $p(y|\boldsymbol{x})$  using  $p(\boldsymbol{x}|y)$ .
- **Assume:**  $p(\boldsymbol{x}|y)p(y) > 0, \forall \boldsymbol{x}, y.$

$$\bullet \frac{p(\boldsymbol{x}|y = 1)p(y = 1)}{p(\boldsymbol{x}|y = 1)p(y = 1) + p(\boldsymbol{x}|y = -1)p(y = -1)} = \frac{1}{1 + \boxed{\frac{p(\boldsymbol{x}|y = -1)p(y = -1)}{p(\boldsymbol{x}|y = 1)p(y = 1)}}}$$

# Modelling $p(y|\mathbf{x})$

---

- We can rewrite  $p(y|\mathbf{x})$  using the ratio  $\frac{p(\mathbf{x}|y = -1)p(y = -1)}{p(\mathbf{x}|y = 1)p(y = 1)}$ .
- $p(y = 1|\mathbf{x}) = \frac{1}{1 + \frac{p(\mathbf{x}|y = -1)p(y = -1)}{p(\mathbf{x}|y = 1)p(y = 1)}}$
- This derivation shows an important difference between generative/discriminative modelling:
  - Generative learning models **class density**  $p(\mathbf{x}|y)$
  - Discriminative learning models **density ratio**  $\frac{p(\mathbf{x}|y = -1)}{p(\mathbf{x}|y = 1)}$ !

# Modelling Density Ratio

---

- Clearly, modelling density ratio  $\frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=-1)}$  requires a whole lot less assumptions on your class densities.
- Models on  $p(\mathbf{x}|y)$   $\Rightarrow$  Models  $\frac{p(\mathbf{x}|y=-1)}{p(\mathbf{x}|y=1)}$
- Models on  $\frac{p(\mathbf{x}|y=-1)}{p(\mathbf{x}|y=1)}$   $\nrightarrow$  Models  $p(\mathbf{x}|y)$

# Modelling Log-Density Ratio

---

$$\bullet p(y = 1 | \mathbf{x}) = \frac{1}{1 + \frac{p(\mathbf{x} | y = -1)p(y = -1)}{p(\mathbf{x} | y = 1)p(y = 1)}}$$
$$\Rightarrow p(y = 1 | \mathbf{x}, \mathbf{w}) := \frac{1}{1 + \exp(-f(\mathbf{x}; \mathbf{w}))}$$

$$\bullet \text{We model log ratio, } \log \frac{p(\mathbf{x} | y = 1)p(y = 1)}{p(\mathbf{x} | y = -1)p(y = -1)} \text{ as } f(\mathbf{x}; \mathbf{w})$$

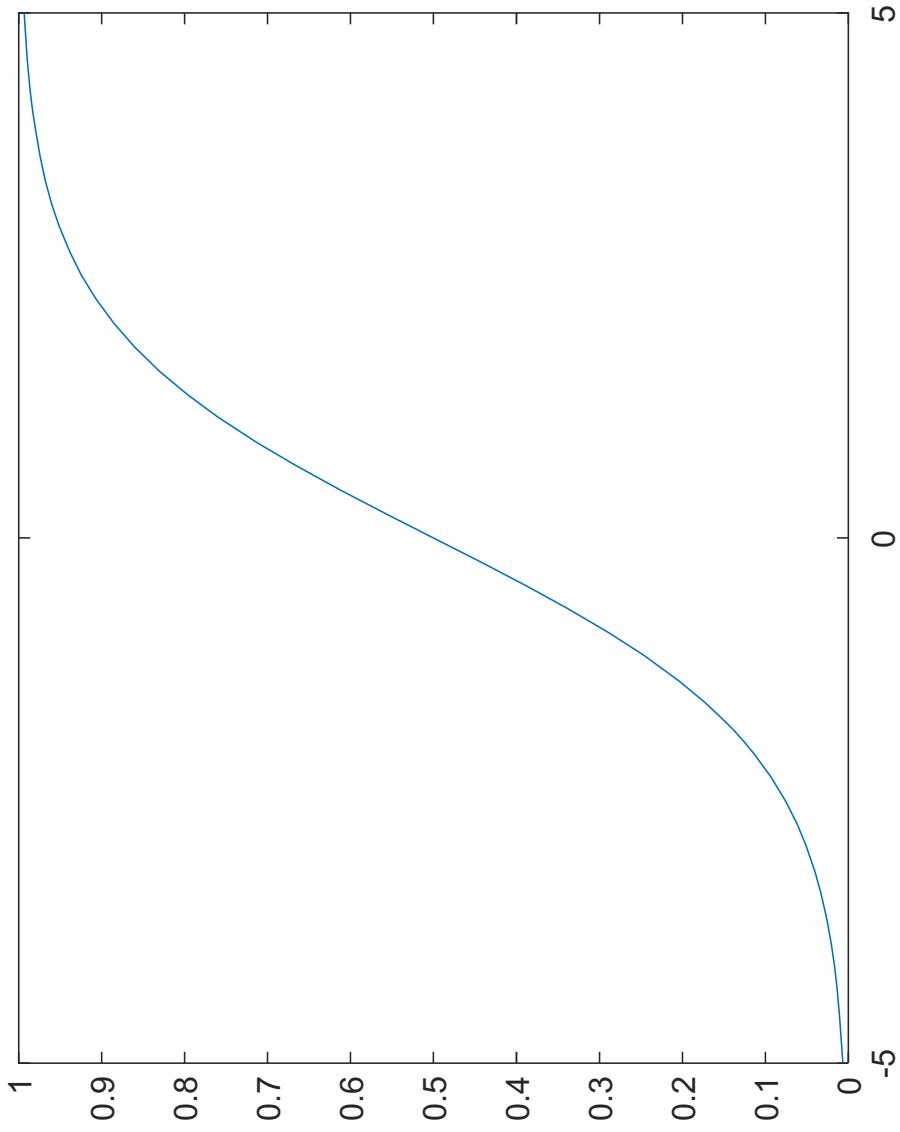
- Like density estimation, it is better to work with log-ratio rather than the ratio itself.

# Generalized Linear Model

---

- As usual,  $f(\mathbf{x}; \mathbf{w}) = \langle \mathbf{w}', \mathbf{x} \rangle + w_0$ .
- Let  $\sigma(t) := \frac{1}{1+\exp(-t)}$ , “sigmoid function”
  - The model for  $p(y|\mathbf{x}; \mathbf{w}) := \sigma(f(\mathbf{x}; \mathbf{w}))$  is merely a linear function wrapped by a non-linear transform.
  - We call  $\sigma(f(\mathbf{x}; \mathbf{w}))$  a “generalized linear model”. This model is widely used in places beyond classification.

Sigmoid Function  $\sigma(t) := \frac{1}{1 + \exp(-t)}$



10

# Modelling Log-Density Ratio

$$\bullet p(y = -1 | \mathbf{x}) = \frac{1}{1 + \frac{p(\mathbf{x} | y = +1)^{p(y=+1)}}{p(\mathbf{x} | y = -1)^{p(y=-1)}}}$$
$$\Rightarrow p(y = -1 | \mathbf{x}, \mathbf{w}) := \frac{1}{1 + \exp(f(\mathbf{x}; \mathbf{w}))}$$

- In  $p(y = -1 | \mathbf{x})$ ,  $\frac{p(\mathbf{x} | y = +1)^{p(y=+1)}}{p(\mathbf{x} | y = -1)^{p(y=-1)}}$  occurs, which is the exact inverse of the ratio appeared in  $p(y = 1 | \mathbf{x})$ . This ratio is modelled by  $\exp(f(\mathbf{x}; \mathbf{w}))$ .

- To simplify our model, let us write

$$\bullet p(y | \mathbf{x}; \mathbf{w}) := \sigma(f(\mathbf{x}; \mathbf{w}) \cdot y)$$

# Estimate $p(y|\mathbf{x}; \mathbf{w})$ from $D$

---

- Assuming the IID-ness on  $D$ .
- Likelihood:  $p(D|\mathbf{w}) = \prod_{i \in D} p(y_i|\mathbf{x}_i; \mathbf{w})$ ,
- Just like what we did for regression tasks.

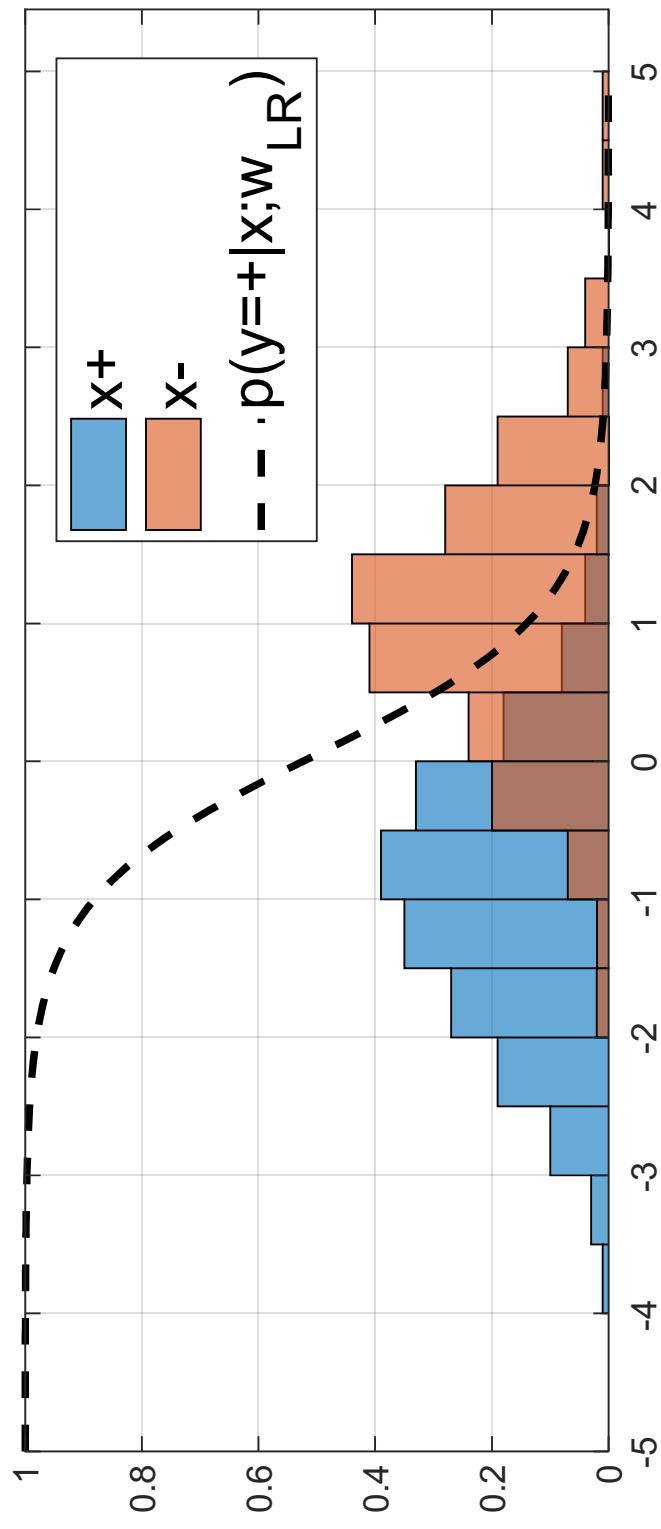
$$\begin{aligned} &\bullet \text{ MLE for } p(y|\mathbf{x}; \mathbf{w}): \\ &\bullet \mathbf{w}_{\text{MLE}} = \operatorname{argmax}_{\mathbf{w}} \log \prod_{i \in D} p(y_i|\mathbf{x}_i; \mathbf{w}) \\ &\quad = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log p(y_i|\mathbf{x}_i; \mathbf{w}) \\ &\quad = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) \end{aligned}$$

# Logistic Regression

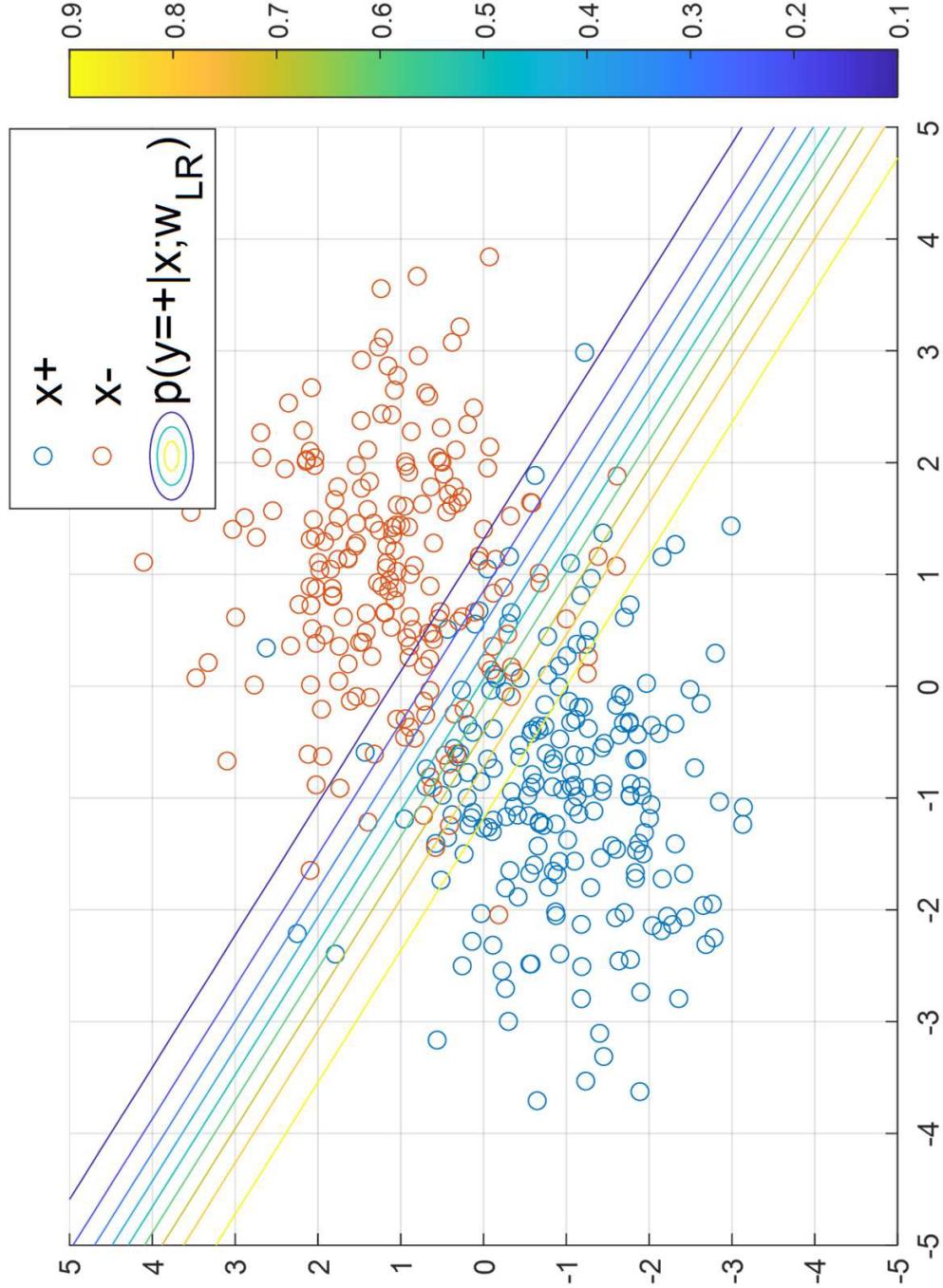
---

- This MLE procedure is also called Logistic Regression.
- **This. Is. Not. A. Regression!**

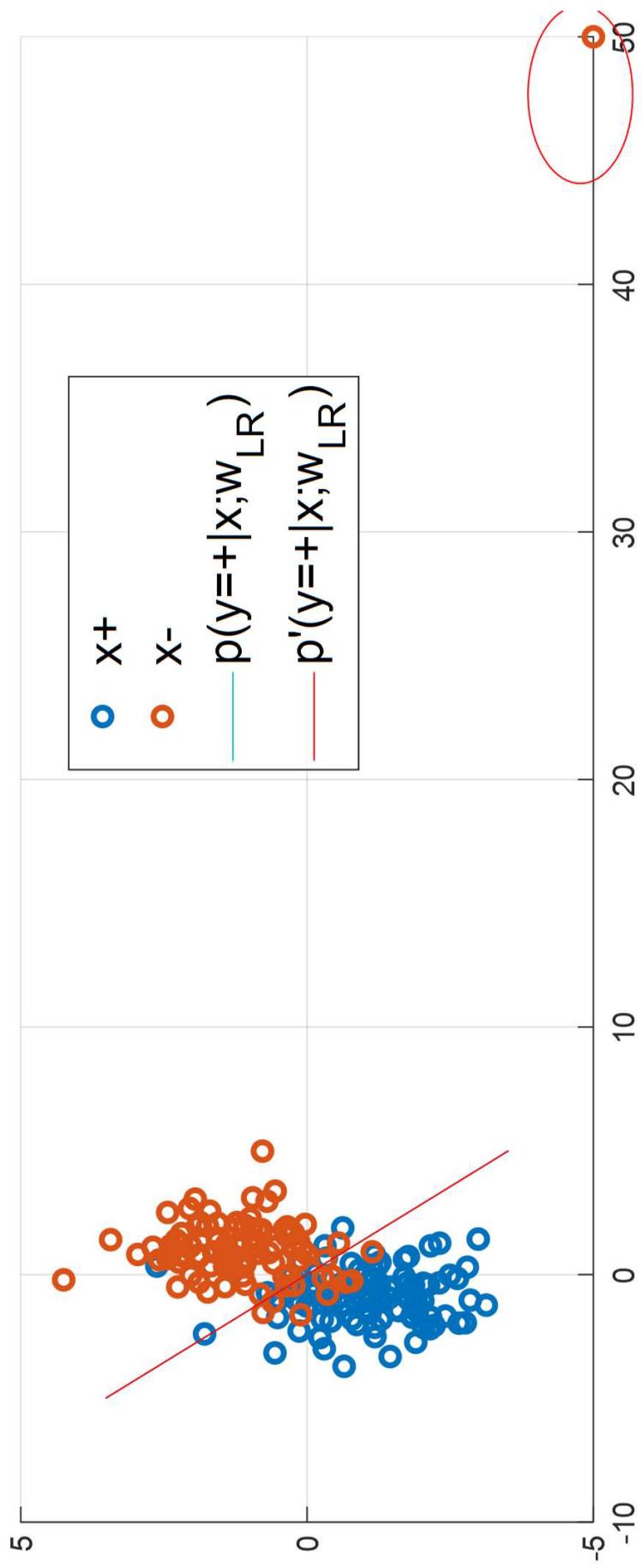
# Logistic Regression



# Logistic Regression 2D

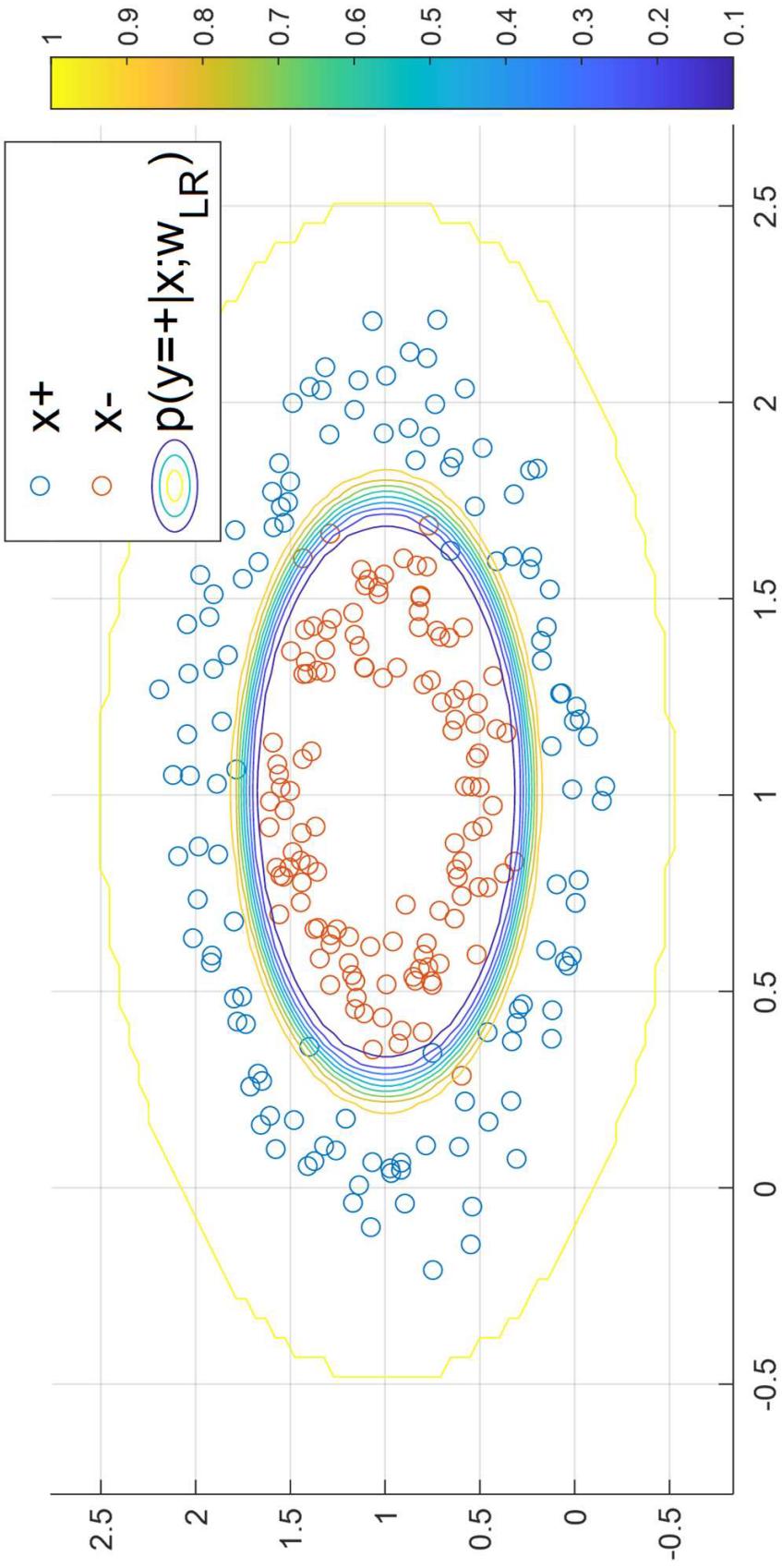


# Robustness of Logistic Regression



Unlike LS classifier, LR is not affected by outliers that are far away from the decision boundary. Why?

# Logistic Regression with Feature Transform $\phi(x)$



- Since  $f(x; w) = \langle w, x \rangle$  still takes a linear form, we can replace  $x$  with  $\phi(x)$  to create a non-linear classifier.
- $\phi$  can be Poly. Trigonometric, or RBF.

# Estimating $p(y|\mathbf{x}; \mathbf{w})$

---

- We can assume priors on  $\mathbf{w}$ , then

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log(\sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i)) p(\mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) + \log p(\mathbf{w}) \end{aligned}$$

- We can also use the full prob. approach

$$\begin{aligned} p(y|\mathbf{x}) &= \int p(y|\mathbf{x}; \mathbf{w}) p(\mathbf{w}|D) d\mathbf{w} \\ &\propto \int p(y|\mathbf{x}; \mathbf{w}) p(D|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} \\ &\propto \int \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) \prod_{i \in D} \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) p(\mathbf{w}) d\mathbf{w} \end{aligned}$$

- Unlike regression using MVN models, we cannot calculate this integral in closed form. See PRML 4.4, 4.5.

# Multi-class Logistic Regression

---

- It is easy to extend logistic regression to a multi-class classification problem.

$$p(y=1|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|y=1)p(y=1)}{\sum_k p(\boldsymbol{x}|y=k)p(y=k)}$$

Marginalization is no longer with respect to a binary  $y$ !

- This expression enables an elegant **expression** of logistic regression **objective** using one-hot encoding.

# One-hot Logistic Regression

---

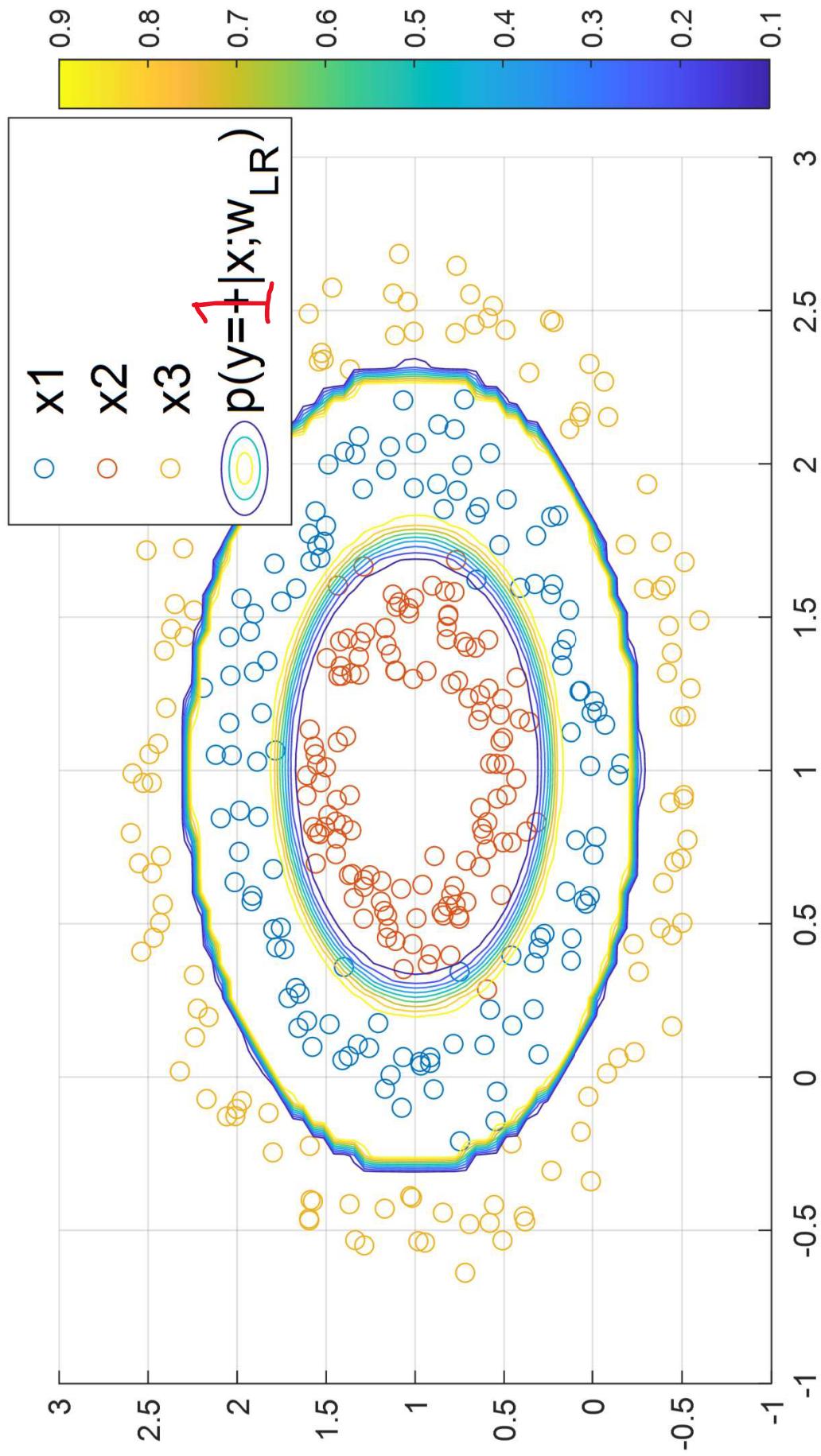
- $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \tilde{\mathbf{x}}$ ,  $\mathbf{w} \in R^{d \times K}$ ,  $\tilde{\mathbf{x}} := [\mathbf{x}^\top, \mathbf{1}]^\top$
- Use “one hot encoding”:  $y_i \in \{1 \dots K\} \Rightarrow \mathbf{t}_i \in R^K$
- $\mathbf{w}_{\text{MLE}} = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}), \mathbf{t}_i)$
- where  $\sigma(f, \mathbf{t}) := \frac{\exp \langle f, \mathbf{t} \rangle}{\sum_k \exp f^{(k)}}$ .
- **Homework:** What is the probabilistic interpretation of  $f$ ?
- If prediction is given by  $\operatorname{argmax}_y p(y | \mathbf{x}; \mathbf{W})$ , it corresponds to multi-class decision rule we saw in previous lecture.  
Why?

# Multi-class Classification

~~Previous Lecture~~

- Rather than relying on sign of  $f$  to make predictions, we estimate  $K$  functions:
- $\{f_k(\mathbf{x}; \mathbf{w}_k)\}_{k=1}^K$
- Given an  $\mathbf{x}$ , prediction is  $\hat{k}$  if  $f_{\hat{k}}(\mathbf{x}; \mathbf{w}_{\hat{k}}) > f_j(\mathbf{x}; \mathbf{w}_j)$ ,  $\forall j$
- **Problem:**  $f_k$  does not have a simple geometry interpretation anymore.
- However,  $f_k$  does have probabilistic interpretation.

# Multi-class Logistic Regression



# Implementation of Logistic Regression

---

- Unlike LS, LR does not have a closed form solution.

- It means, to find  $\mathbf{w}_{\text{MLE}}$ , we need to solve

$$\operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i)$$

- numerically!!

- The implementation of this algorithm requires some knowledge on numerical optimization, which is not introduced in this class.

- Fortunately, numerical optimization packages are readily available in many programming languages.

# Conclusion

---

- Discriminative classification models **density ratio** while generative classification models **class densities**.
- When log-ratio is modelled by  $f(x; w) := \langle w', x \rangle + w_0$ , the model for the class posterior  $p(y|x)$  is called **generalized linear model**.
- The MLE solution for generalized linear model is called **logistic regression**.
  - whose solution requires numerical optimization.

# Homework

---

- What are the **decision functions** given by a binary logistic regression? (hint:  $p(y|x; w) - .5$  is one of them)
- Prove: if  $p(x|y = 1)$  and  $p(x|y = -1)$  are MVN with shared covariance matrix  $\Sigma$  but different means  $\mu_+, \mu_-$ .
  - 1.  $\exists w^*$  such that  $p(y|x) = \sigma(\langle x; w'^* \rangle + w_0'^*)y)$
  - 2. find  $w^*$
- Show the probabilistic interpretation of multiclass logistic regression

# Jensen Shannon Divergence (Challenging)

- Similar to KL divergence, Jensen Shannon divergence is a discrepancy measure between two probability density functions  $p$  and  $q$ .

$$\bullet JS[p, q] := \frac{1}{2} E_p \left[ \log \frac{p(x)}{.5p(x) + .5q(x)} \right] + \frac{1}{2} E_q \left[ \log \frac{q(x)}{.5p(x) + .5q(x)} \right].$$

- How is the LR objective related to  $JS[p, q]$  when  $p(y=1) = p(y=-1)$ ?

- Hint: What is the maximiser of the following problem?

- $\arg\max_t E_p[\log t(x)] + E_q[\log(1 - t(x))]$ , where  $t$  is a function  $t: R^d \rightarrow R$ ,  $t \in (0, 1)$ .