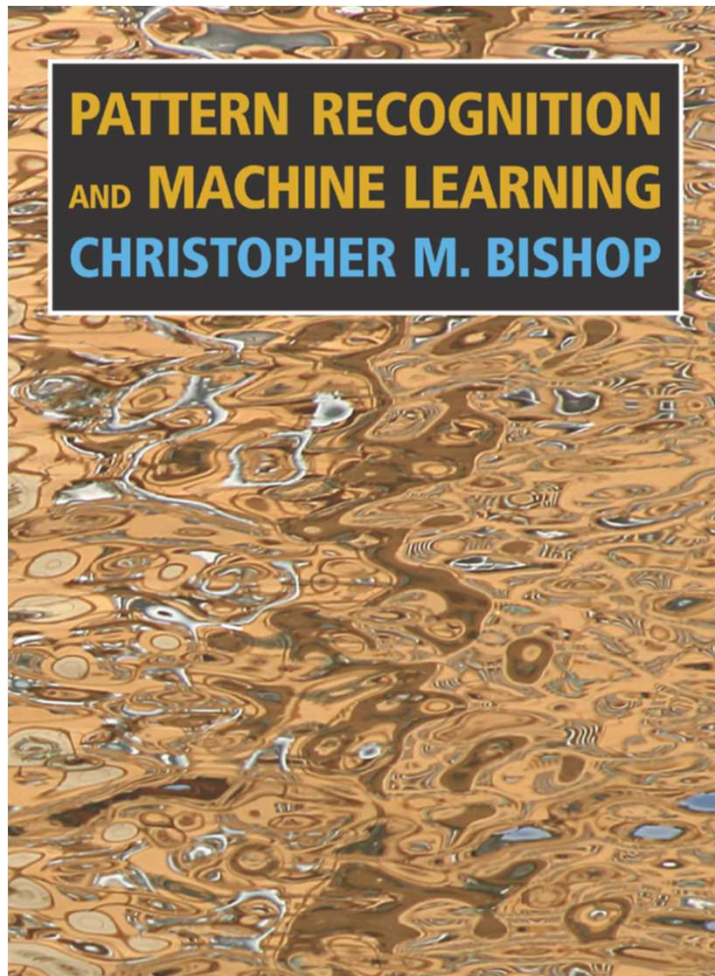# Discriminative Classifiers

Song Liu (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 $p(y|\boldsymbol{x})$ given dataset $D$.

- **Step 1.** Making a model assumption $p(y|\boldsymbol{x}; \boldsymbol{w})$.
- **Step 2.** Construct the likelihood function $p(D|\boldsymbol{w})$.
- **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:

- $p(y|\boldsymbol{x}) = \dfrac{p(\boldsymbol{x}|y)p(y)}{p(\boldsymbol{x})} = \dfrac{p(\boldsymbol{x}|y)p(y)}{\sum_{y'} p(x,y')} = \dfrac{p(\boldsymbol{x}|y)p(y)}{\sum_{y'} p(\boldsymbol{x}|y')p(y')}$ so

  <span style="color:red">Marginalization!</span>

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

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

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

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

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

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

- $\dfrac{p(\boldsymbol{x}|y = 1)p(y=1)}{p(\boldsymbol{x}|y = 1)p(y=1)+p(\boldsymbol{x}|y = -1)p(y=-1)} = \dfrac{1}{1+\dfrac{p(\boldsymbol{x}|y = -1)p(y=-1)}{p(\boldsymbol{x}|y = 1)p(y=1)}}$

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

- We can rewrite $p(y|\boldsymbol{x})$ using the ratio $\dfrac{p(\boldsymbol{x}|y=-1)p(y=-1)}{p(\boldsymbol{x}|y=1)p(y=1)}$:

- $p(y=1|\boldsymbol{x}) = \dfrac{1}{1+\dfrac{p(\boldsymbol{x}|y=-1)p(y=-1)}{p(\boldsymbol{x}|y=1)p(y=1)}}$

- This derivation shows an important difference between generative/discriminative modelling:

- Generative learning models class density $p(\boldsymbol{x}|y)$
- Discriminative learning models density ratio $\dfrac{p(\boldsymbol{x}|y=-1)}{p(\boldsymbol{x}|y=1)}$!

6

# Modelling Density Ratio

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

- Models on $p(\boldsymbol{x}|y) \Rightarrow$ Models $\dfrac{p(\boldsymbol{x}|y=-1)}{p(\boldsymbol{x}|y=1)}$

- Models on $\dfrac{p(\boldsymbol{x}|y=-1)}{p(\boldsymbol{x}|y=1)} \not\Rightarrow$ Models $p(\boldsymbol{x}|y)$

# Modelling Log-Density Ratio

- $p(y = 1|x) = \dfrac{1}{1+\dfrac{p(x|y = -1)p(y=-1)}{p(x|y = 1)p(y=1)}}$

  $\Rightarrow p(y = 1|x, w) := \dfrac{1}{1 + \exp(-f(x; w))}$
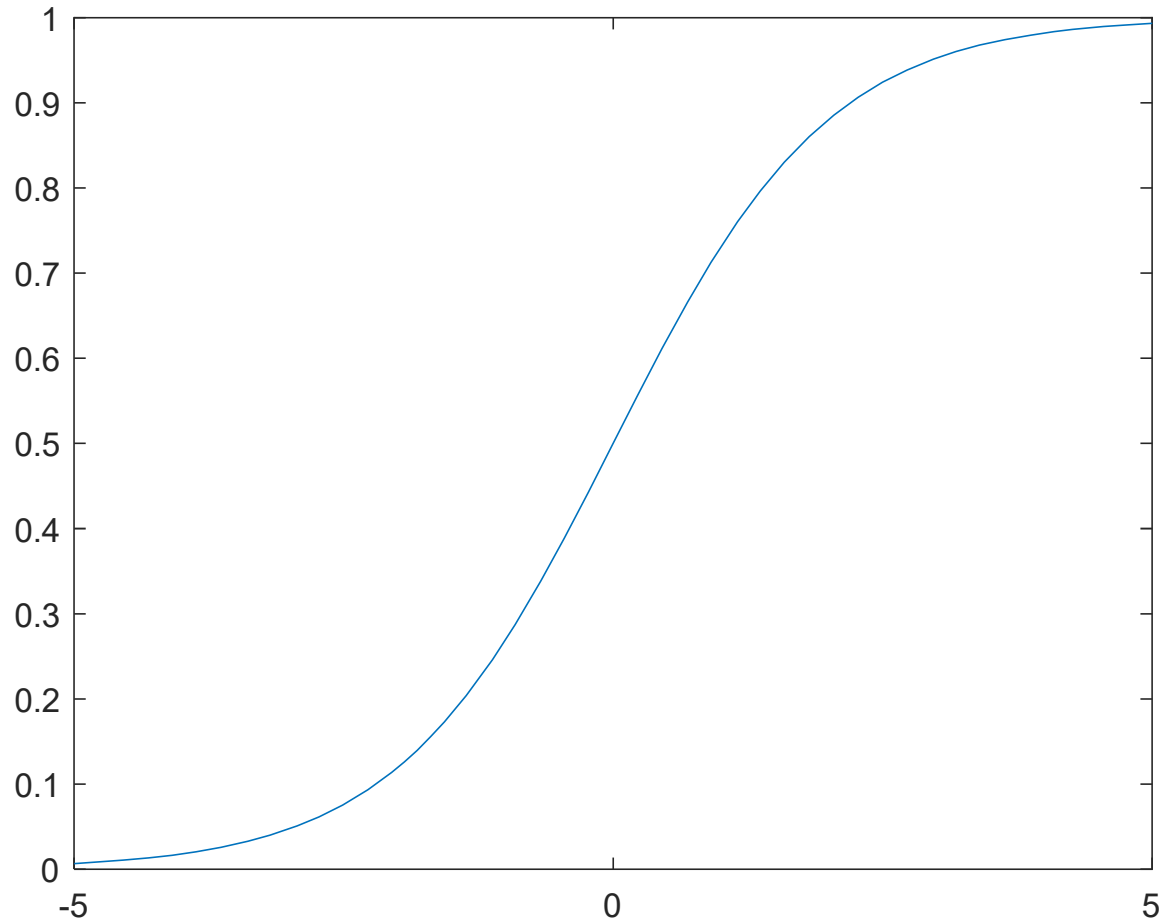
- We model log ratio, $\log \dfrac{p(x|y = 1)p(y=1)}{p(x|y = -1)p(y=-1)}$ as $f(x; w)$

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

# Generalized Linear Model

- As usual, $f(x; w) = \langle w', x \rangle + w_0$.

- Let $\sigma(t) := \dfrac{1}{1+\exp(-t)}$ , "sigmoid function"

- The model for $p(y|x; w) := \sigma\big(f(x; w)\big)$ is merely a linear function wrapped by a non-linear transform.

- We call $\sigma\big(f(x; w)\big)$ a "generalized linear model". This model is widely used in places beyond classification.

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

# Modelling Log-Density Ratio

- $p(y = -1|\boldsymbol{x}) = \dfrac{1}{1+\dfrac{p(\boldsymbol{x}|y = +1)p(y=+1)}{p(\boldsymbol{x}|y = -1)p(y=-1)}}$

$\Rightarrow p(y = -1|\boldsymbol{x}, \boldsymbol{w}) := \dfrac{1}{1 + \exp\big(f(\boldsymbol{x}; \boldsymbol{w})\big)}$

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

- To simplify our model, let us write
- $p(y|\boldsymbol{x}; \boldsymbol{w}) := \sigma(f(\boldsymbol{x}; \boldsymbol{w}) \cdot y)$

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

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

- MLE for $p(y|\boldsymbol{x}; \boldsymbol{w})$:
- $\boldsymbol{w}_{\mathrm{MLE}} = \mathrm{argmax}_{\boldsymbol{w}} \log \prod_{i \in D} p(y_i|\boldsymbol{x}_i; \boldsymbol{w})$

$$= \mathrm{argmax}_{\boldsymbol{w}} \sum_{i \in D} \log p(y_i|\boldsymbol{x}_i; \boldsymbol{w})$$

$$= \color{red}{\mathrm{argmax}_{\boldsymbol{w}} \sum_{i \in D} \log \sigma(f(\boldsymbol{x}_i; \boldsymbol{w}) \cdot y_i)}$$
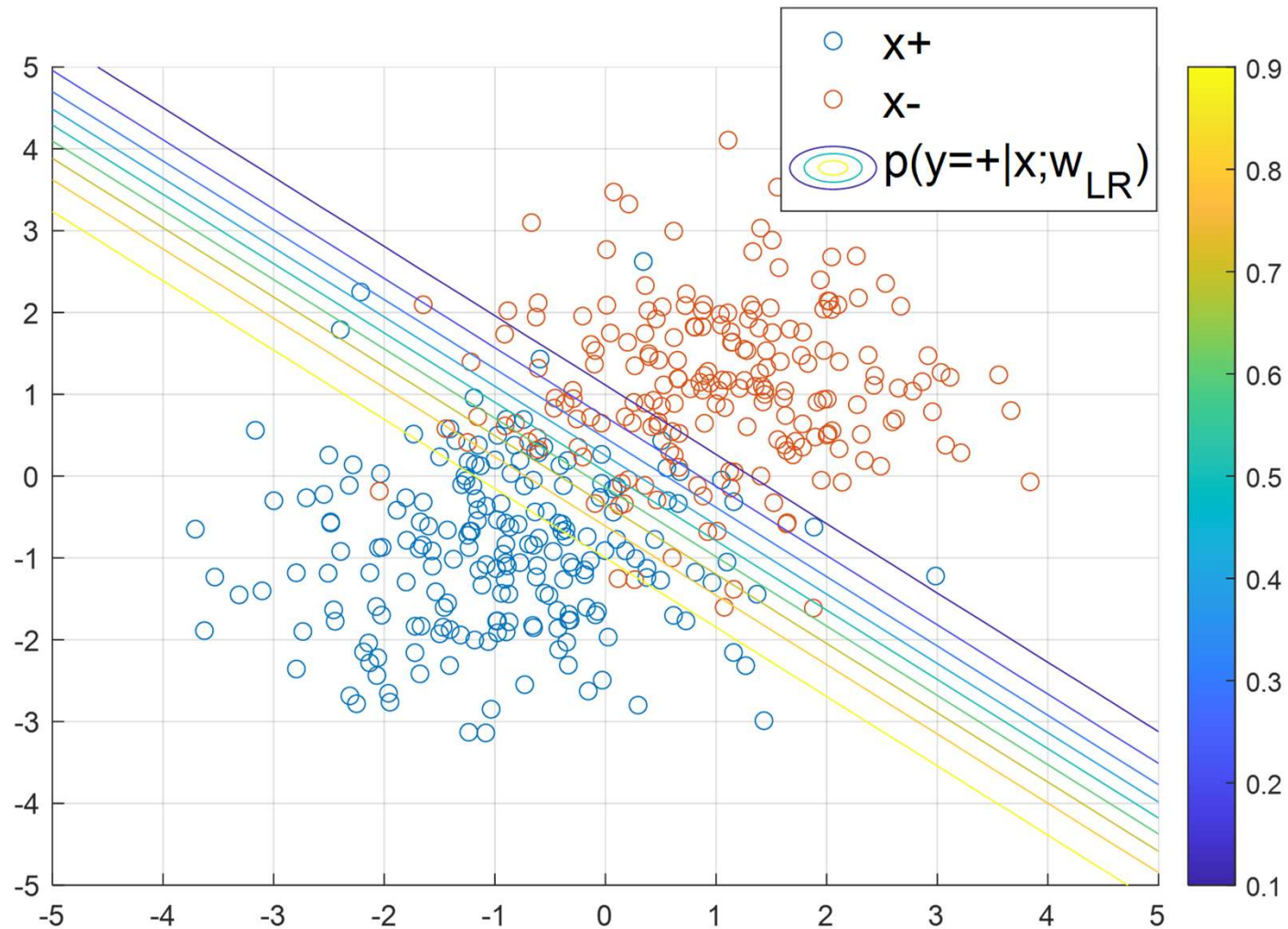
# Logistic Regression

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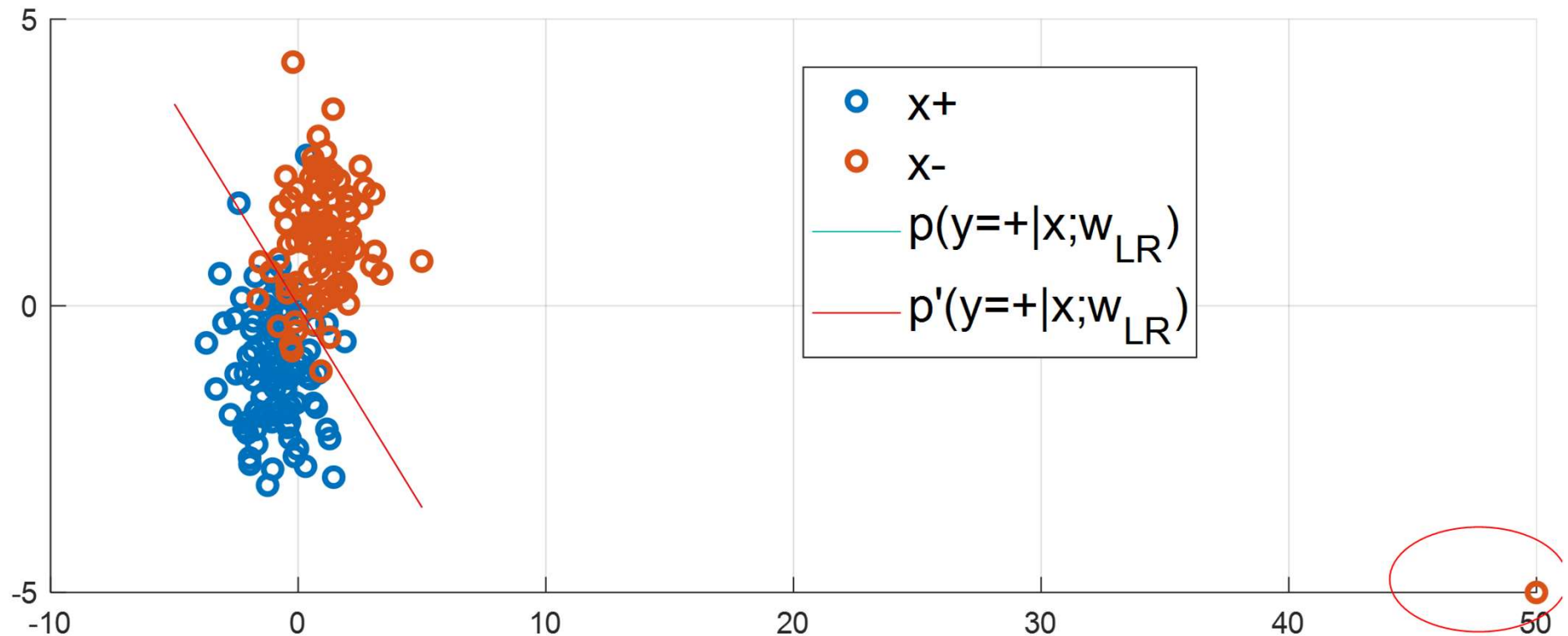- **This. Is. Not. A. Regression!**

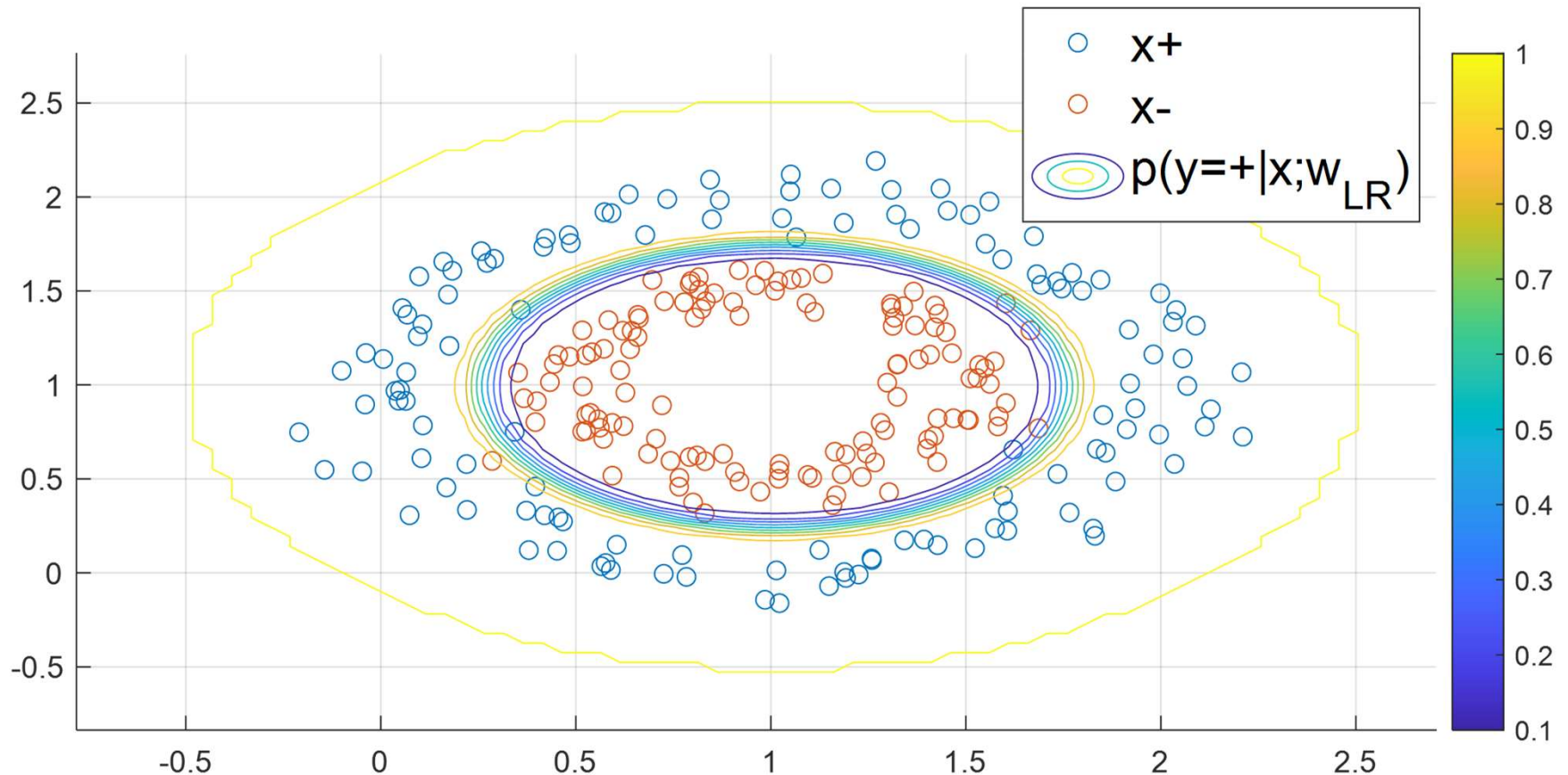# Logistic Regression

# 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. Trignometric, or RBF.

# Estimating $p(y|x; w)$

- We can assume priors on $w$, then
- $w_{\text{MAP}} = \text{argmax}_w \sum_{i \in D} \log(\sigma(f(x_i; w) \cdot y_i) \, p(w))$

$$= \text{argmax}_w \sum_{i \in D} \log \sigma(f(x_i; w) \cdot y_i) + \log p(w)$$

- We can also use the full prob. approach
- $p(y|x) = \int p(y|x; w) p(w|D) dw$

$$\propto \int p(y|x; w) p(D|w) p(w) dw$$

$$\propto \int \sigma(f(x_i; w) \cdot y_i) \prod_{i \in D} \sigma(f(x_i; w) \cdot y_i) \, p(w) dw$$

- 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.

- $p(y = 1|\boldsymbol{x}) = \dfrac{p(\boldsymbol{x}|y = 1)p(y=1)}{\boxed{\sum_k p(\boldsymbol{x}|y = k)p(y=k)}}$

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

- This expression enables an elegant expression of logistic regression objective using one-hot encoding.
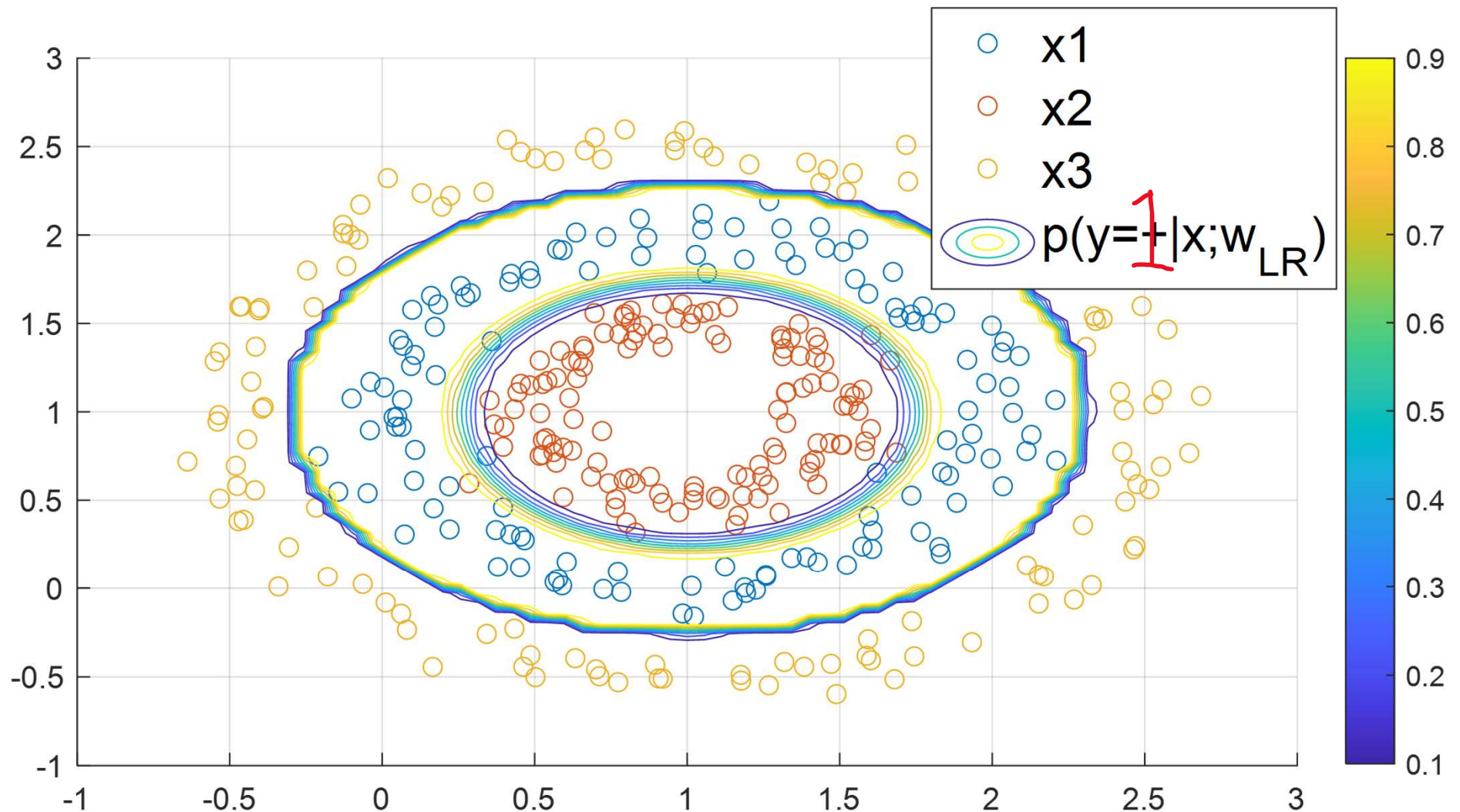
# One-hot Logistic Regression

- $f(x; w) = W^\top \widetilde{x}, W \in R^{d \times K}, \widetilde{x} := [x^\top, 1]^\top$
- Use "**one hot encoding**": $y_i \in \{1 \dots K\} \Rightarrow t_i \in R^K$

- $w_{\text{MLE}} = \text{argmax}_w \sum_{i \in D} \log \sigma(f(x_i; w), t_i)$
- where $\sigma(f, t) := \frac{\exp\langle f, t \rangle}{\sum_k \exp f^{(k)}}$.

- **Homework: What is the probabilistic interpretation of $f$?**
- If prediction is given by $\text{argmax}_y p(y|x; W)$, it corresponds to multi-class decision rule we saw in previous lecture. Why?

20

# Multi-class Classification

- Rather than relying on sign of $f$ to make predictions, we estimate $K$ functions:

- $\{f_k(\boldsymbol{x}; \boldsymbol{w}_k)\}_{k=1}^{K}$

- Given an $\boldsymbol{x}$, prediction is $\hat{k}$ if $f_{\hat{k}}(\boldsymbol{x}; \boldsymbol{w}_{\hat{k}}) > f_j(\boldsymbol{x}; \boldsymbol{w}_j), \forall j$

- **Problem**: $f_k$ does not have a simple geometry interpretation anymore.

- However, $f_k$ does have probabilistic interpretation.

Previous Lecture

# Multi-class Logistic Regression

# Implementation of Logistic Regression

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

- It means, to find $\boldsymbol{w}_{\mathrm{MLE}}$, we need to solve
$$\mathrm{argmax}_{\boldsymbol{w}} \sum_{i \in D} \log \sigma(f(\boldsymbol{x}_i; \boldsymbol{w}) \cdot y_i)$$

- 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 $f(\boldsymbol{x}; \boldsymbol{w}) := \langle \boldsymbol{w}', \boldsymbol{x} \rangle + w_0$, 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 $p(y|x) = \sigma\big((\langle x; w'^*\rangle + w_0'^*)y\big)$
- 2. find $w^*$

- Show the probabilistic interpretation of multiclass logistic regression

# Jensen Shannon Divergence (Challenging)

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

- $\text{JS}[p, q] := \frac{1}{2} \mathrm{E}_p \left[ \log \frac{p(x)}{.5p(x) + .5q(x)} \right] + \frac{1}{2} \mathrm{E}_q \left[ \log \frac{q(x)}{.5p(x) + .5q(x)} \right]$.

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

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

- $argmax_t \mathrm{E}_p[\log t(x)] + \mathrm{E}_q[\log(1 - t(x))]$, where $t$ is a function $t: R^d \to R, t \in (0,1)$.