# Linear Classifiers

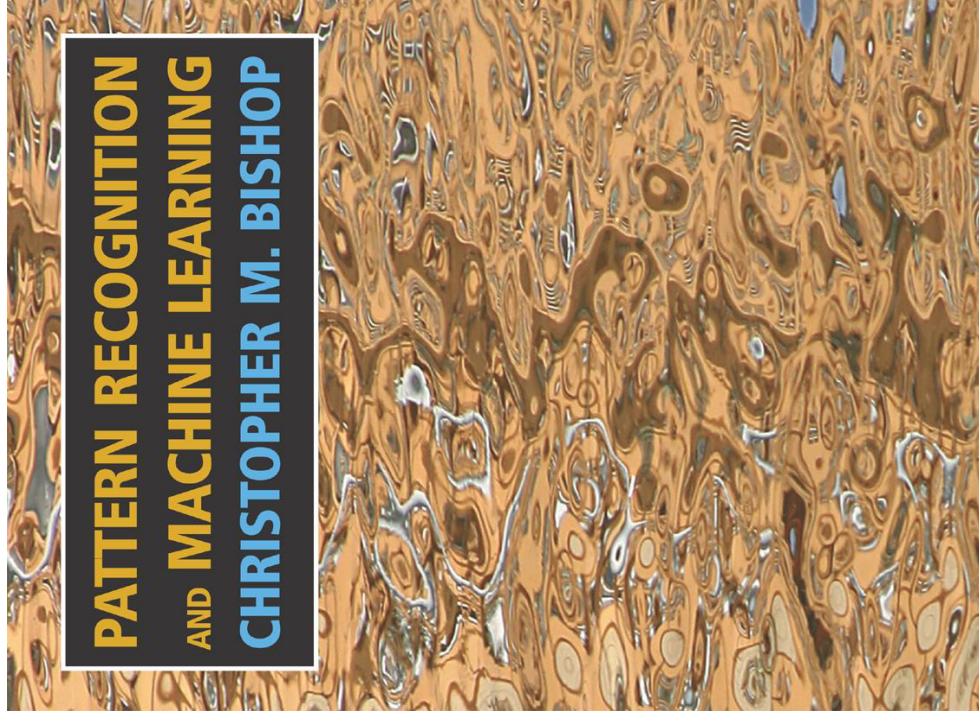Song Liu ([song.liu@bristol.ac.uk](song.liu@bristol.ac.uk))

# Reference

Today's class *roughly* follows Chapter 4-4.2.

Pattern Recognition and Machine Learning

Christopher Bishop, 2006
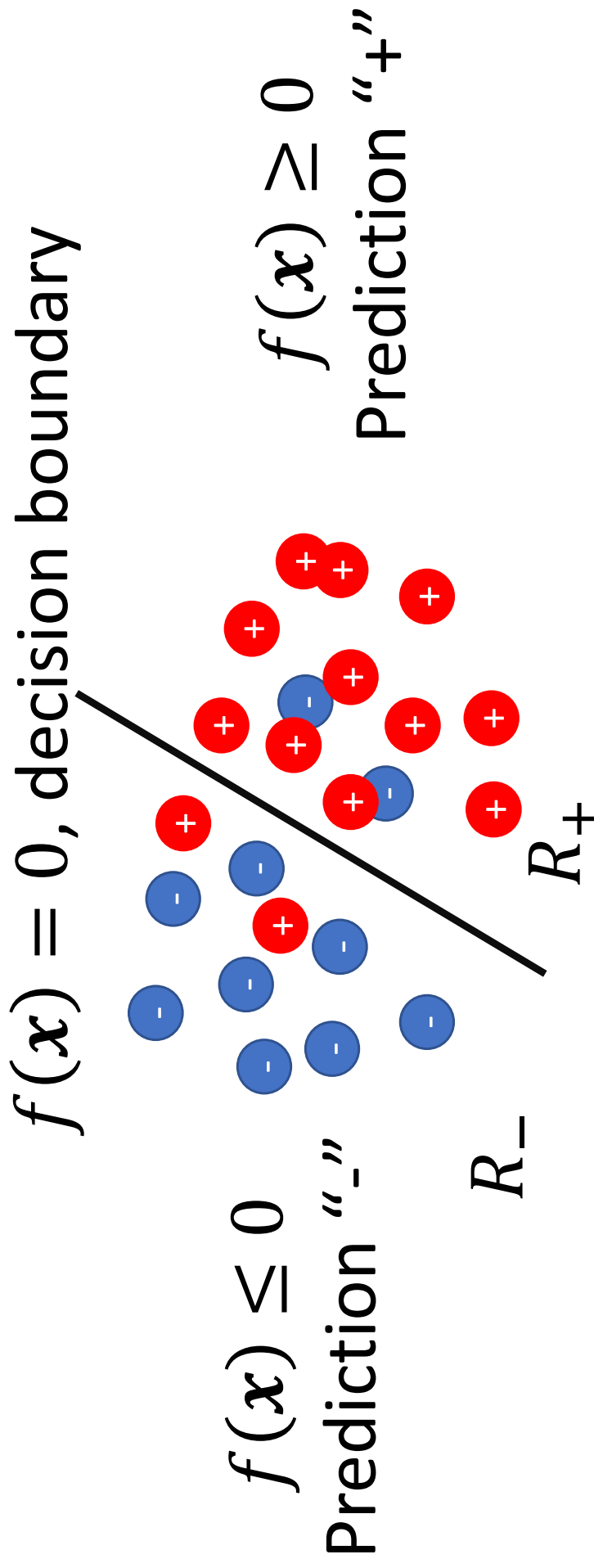


PATTERN RECOGNITION AND MACHINE LEARNING
CHRISTOPHER M. BISHOP

# Outline

- Geometry of decision function

- **Non probabilistic classifiers**

  - Least square classifier
  - Fisher discriminant analysis

- **Probabilistic classifiers**
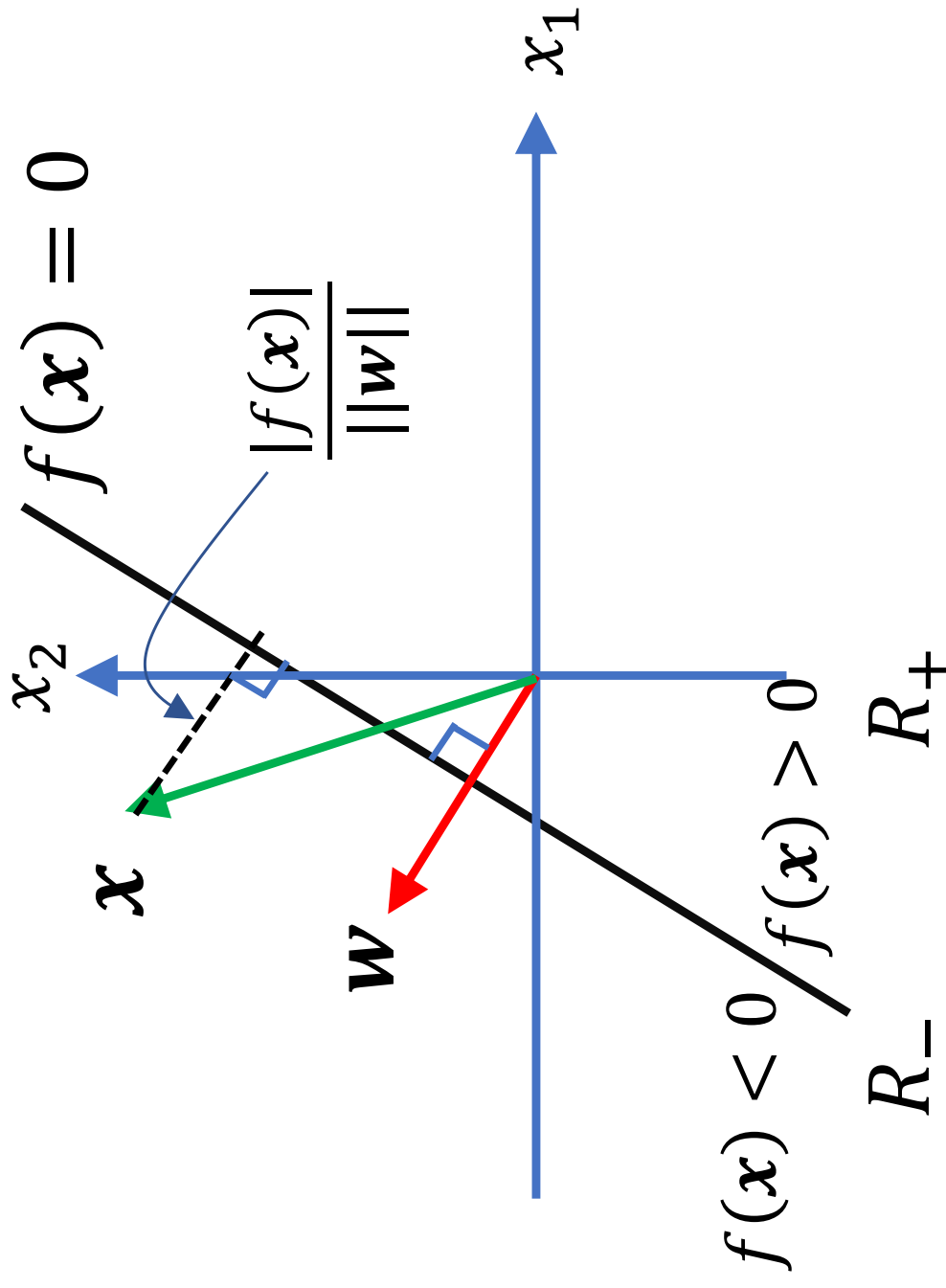
  - Generative Classifiers

# Binary Classification

- **Input:** $x \in R^d$
- **Output:** $y \in \{-1, +1\}$
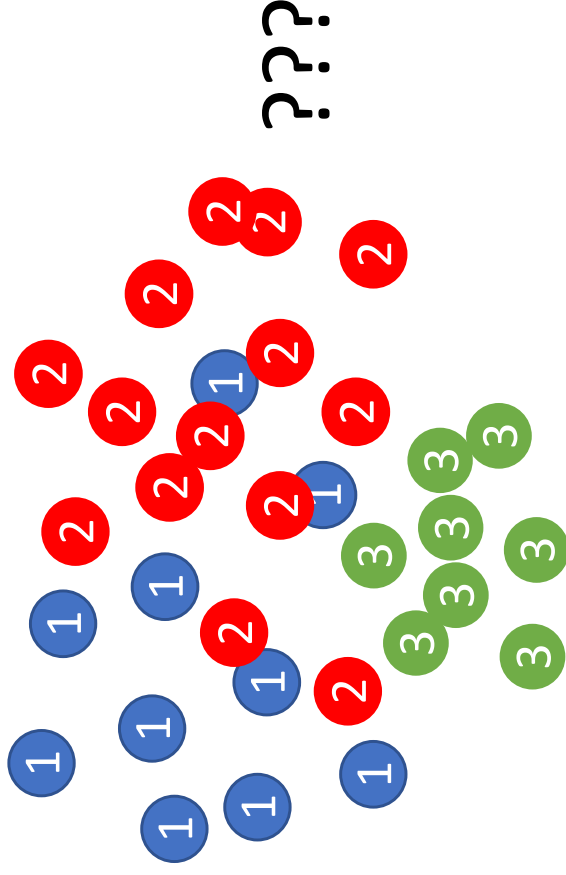- A decision boundary is defined by a function $f(x)$

$f(x) = 0$, decision boundary

$f(x) \leq 0$
Prediction "−"

$R_-$

$f(x) \geq 0$
Prediction "+"

$R_+$

# Geometry of Binary Classification

Suppose $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + w_0$



$$f(\mathbf{x}) = 0$$

$$\frac{|f(\mathbf{x})|}{\|\mathbf{w}\|}$$

$x_2$

$x_1$

$\mathbf{x}$

$\mathbf{w}$

$R_-$

$R_+$

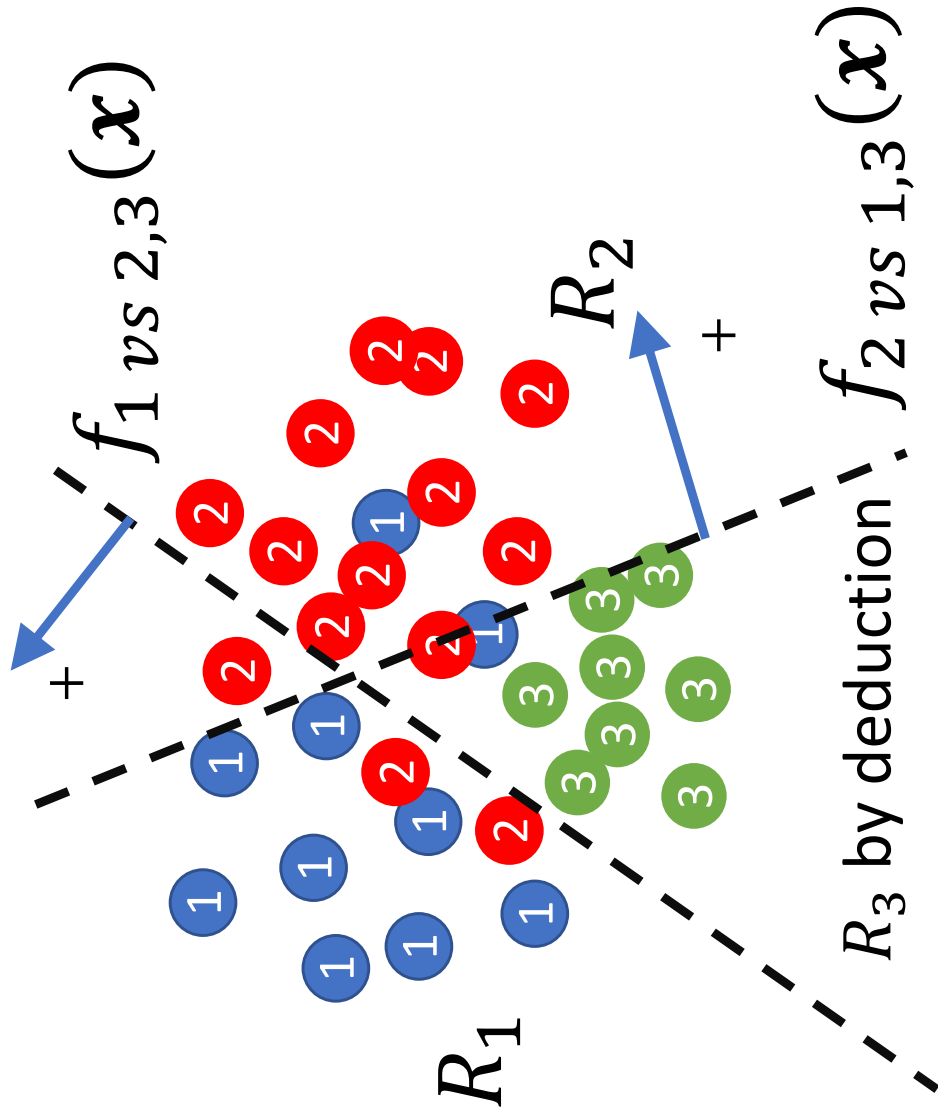$f(\mathbf{x}) < 0$   $f(\mathbf{x}) > 0$

# Multi-class Classification

- **Input:** $x \in R^d$

- **Output:** $y \in \{1 \dots K\}$

- The geometry gets a lot more complicated...

  - Cannot simply check the sign of a single $f(x)$.

???

# One versus The Other

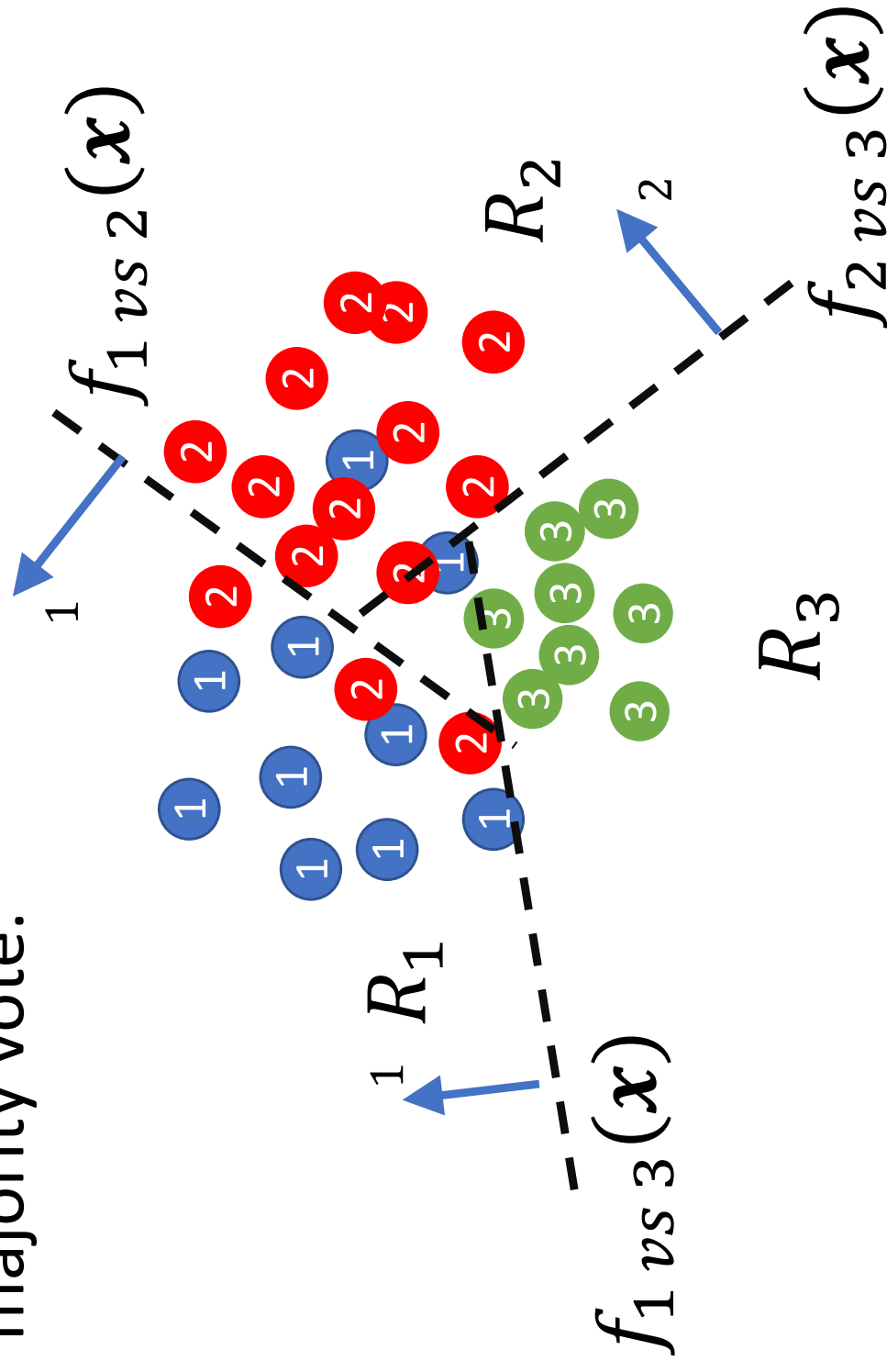- Construct $K - 1$ binary classifiers
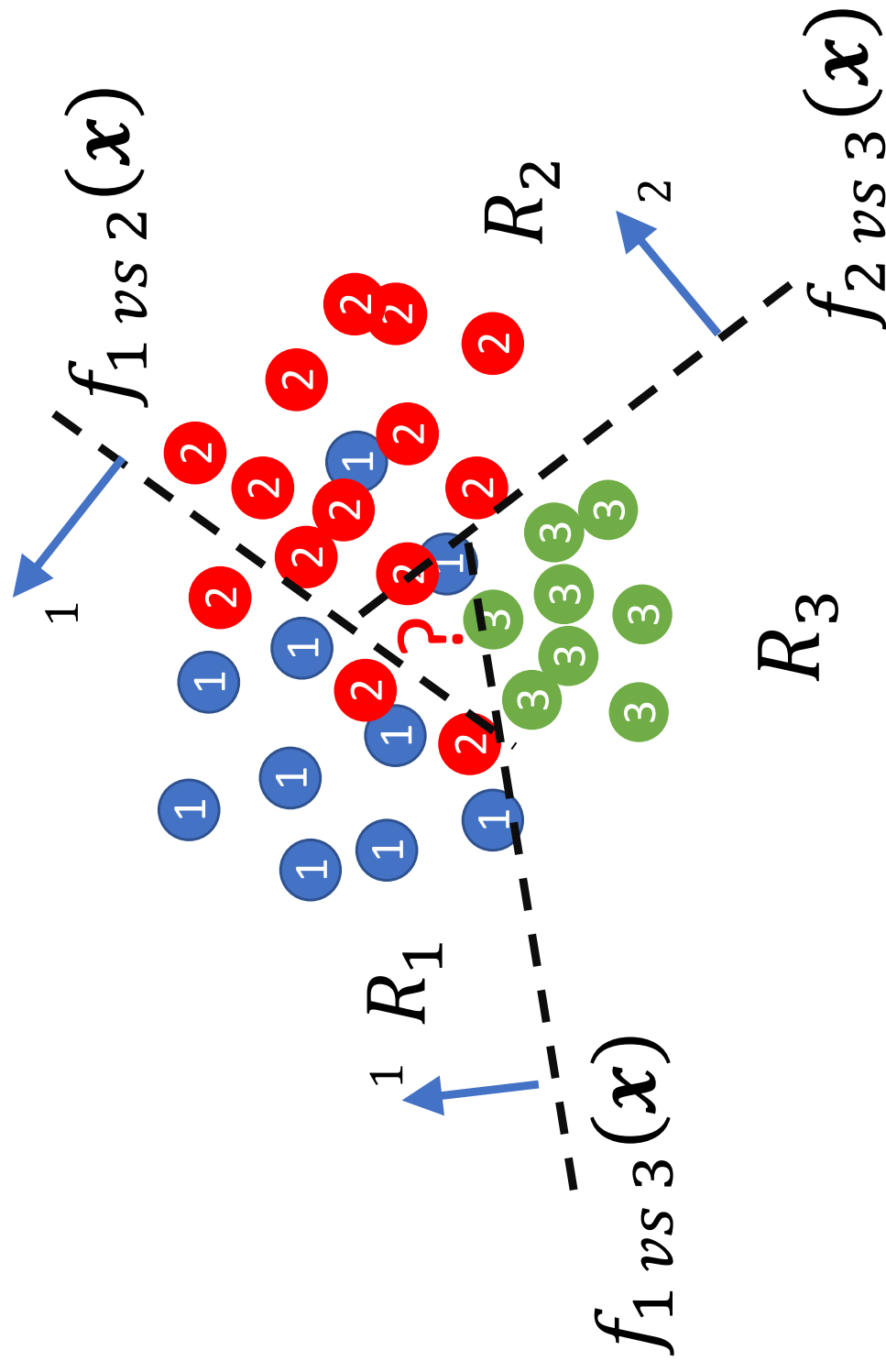- Classify Class $k$ vs. the rest of classes

# One versus The Other

- One versus the other also creates confusion!



$f_{1 \, vs \, 2,3}(x)$

$f_{2 \, vs \, 1,3}(x)$

????

$R_1$

$R_2$

$R_3$ by deduction

# One versus One

- We can create pairwise binary classifiers and check majority vote.



$f_{1\,vs\,2}(\boldsymbol{x})$

$f_{2\,vs\,3}(\boldsymbol{x})$

$f_{1\,vs\,3}(\boldsymbol{x})$

$R_1$

$R_2$

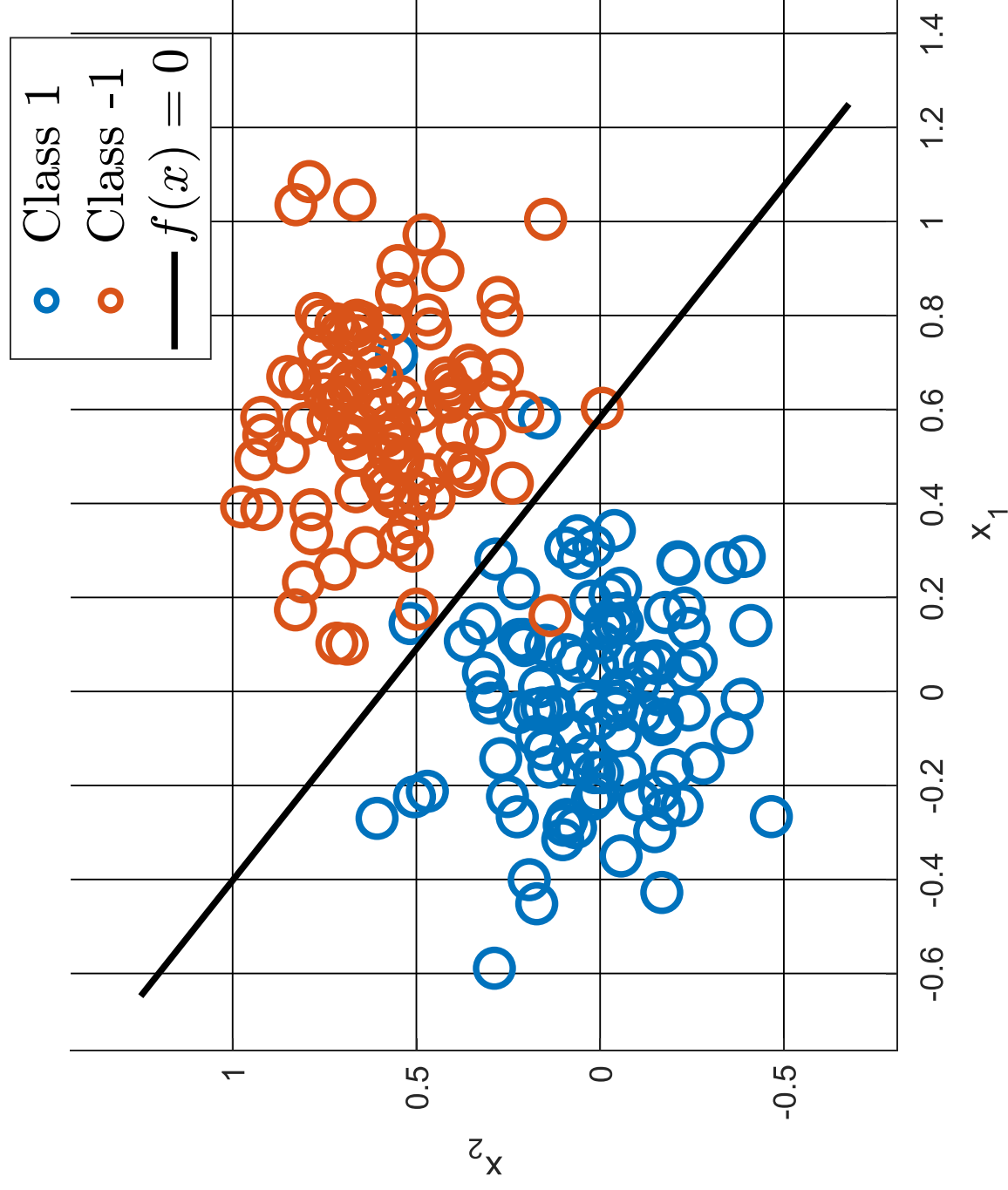$R_3$

# One versus One

- One versus one creates confusion as well...

# Multi-class Classification

- Or...

- rather than relying on sign of $f$ to make predictions, we can fit a vector valued function $\boldsymbol{f}: R^d \to R^K$:

- Given an $\boldsymbol{x}$, prediction is $\hat{k} = \underset{k}{\mathrm{argmax}}\, f^{(k)}(\boldsymbol{x})$,

  - The classification does not have a simple geometry interpretation anymore.

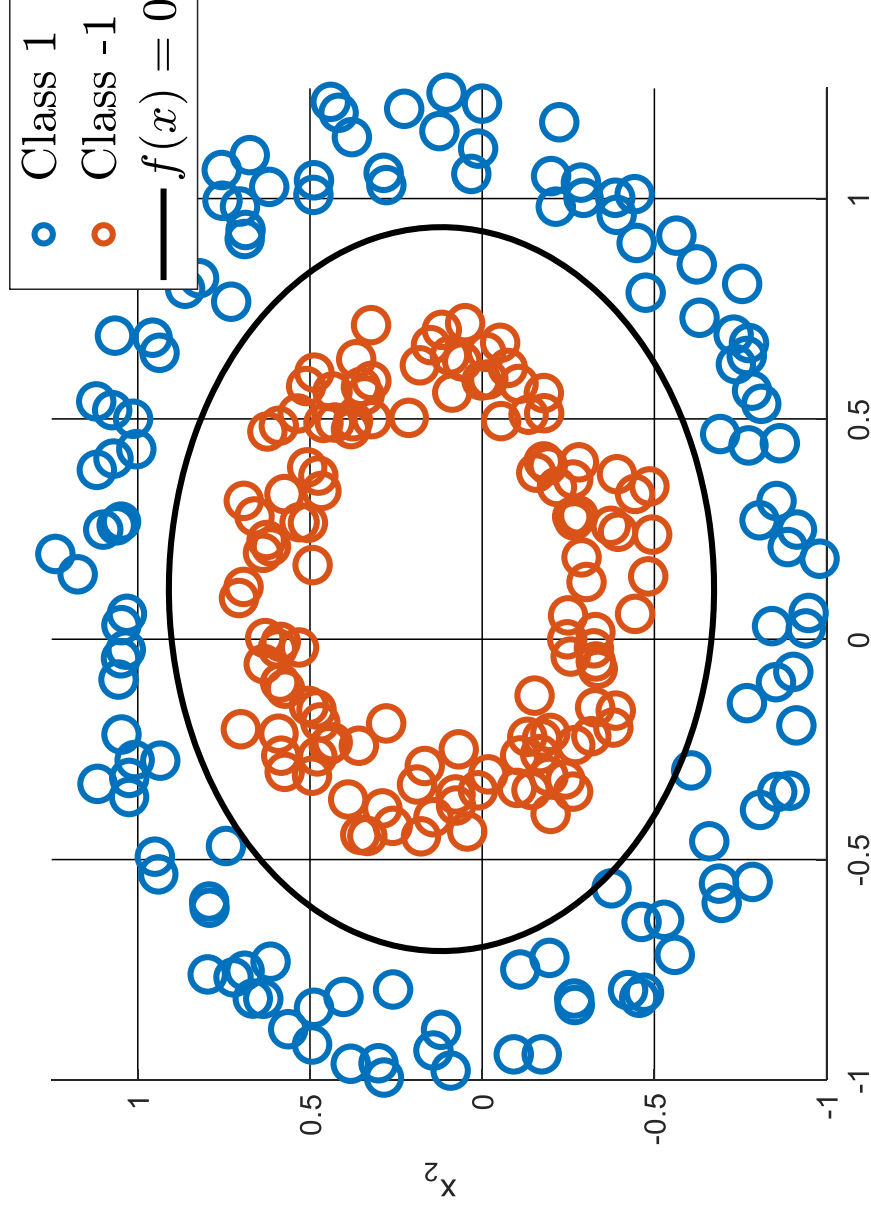  - We will see an example soon.

# Least Squares Classifier

- For binary classification, perform LS on $D$.

- $w_{\text{LS}} := \underset{w}{\arg\min} \sum_{i \in D} [y_i - f(x_i; w)]^2$

  - Now $y_i$ takes binary value 1 or $-1$

- Prediction function $f(x_i; w_{\text{LS}})$.

- The predicted label $\hat{y} := \text{sign}(f(x_i; w_{\text{LS}}))$

# Least Square Classifier

# Least Square Classifier

- You can use feature transform $\boldsymbol{\phi}$ for $f$ as well.

- $f(\boldsymbol{x}; \boldsymbol{w}) := \langle \boldsymbol{w}, \boldsymbol{\phi}(\boldsymbol{x}) \rangle$,

- e.g. poly., trigonometric, RBF, kernel.

# Least Square Classifier

Data may not be separable in the original space but can be separable in the **feature space** created by $\phi$!
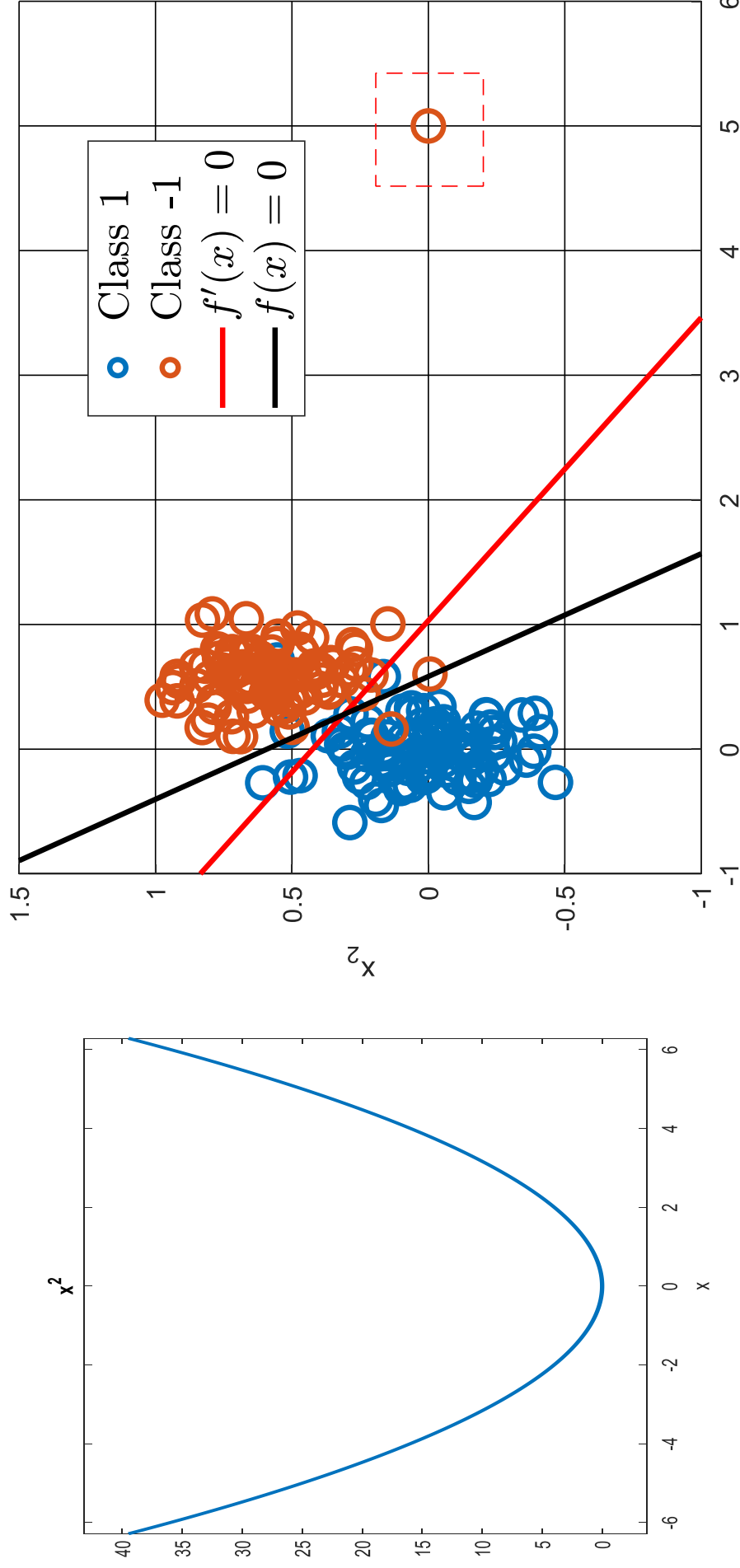
# Multi-class LS classification

- LS can be adapted to multi-class classification.

- Suppose output $y \in \{1 \dots K\}$

- Replace $y_i = k$ in $D$ with $t_i \in \{0,1\}^K$.
  - $t_i^{(k)} = 1$.
  - $t_i^{(j)} = 0, \forall j \neq k$

- **"One-hot encoding"**

- $W_{\mathrm{LS}} := \mathrm{argmin}_W \sum_{i \in D} \|t_i - f(x_i; W)\|^2$

- $W \in R^{(d+1) \times K}, \tilde{x}_i := [x_i^\top, 1]^\top \in R^d, f(x; W) = W^\top \tilde{x}$

- Prediction: $\hat{k} = \mathrm{arg} \max_k f^{(k)}(x; W) = \mathrm{arg} \max_k \left(w_{\mathrm{LS}}^{(k)}\right)^\top \tilde{x}$
  - where $w^{(k)}$ is the $k$-th column of $W$.

# Why not to use LS Classifier?

- Square loss does not make sense in classification tasks.

- Data point far away from decision boundary can influence the decision boundary by a lot.
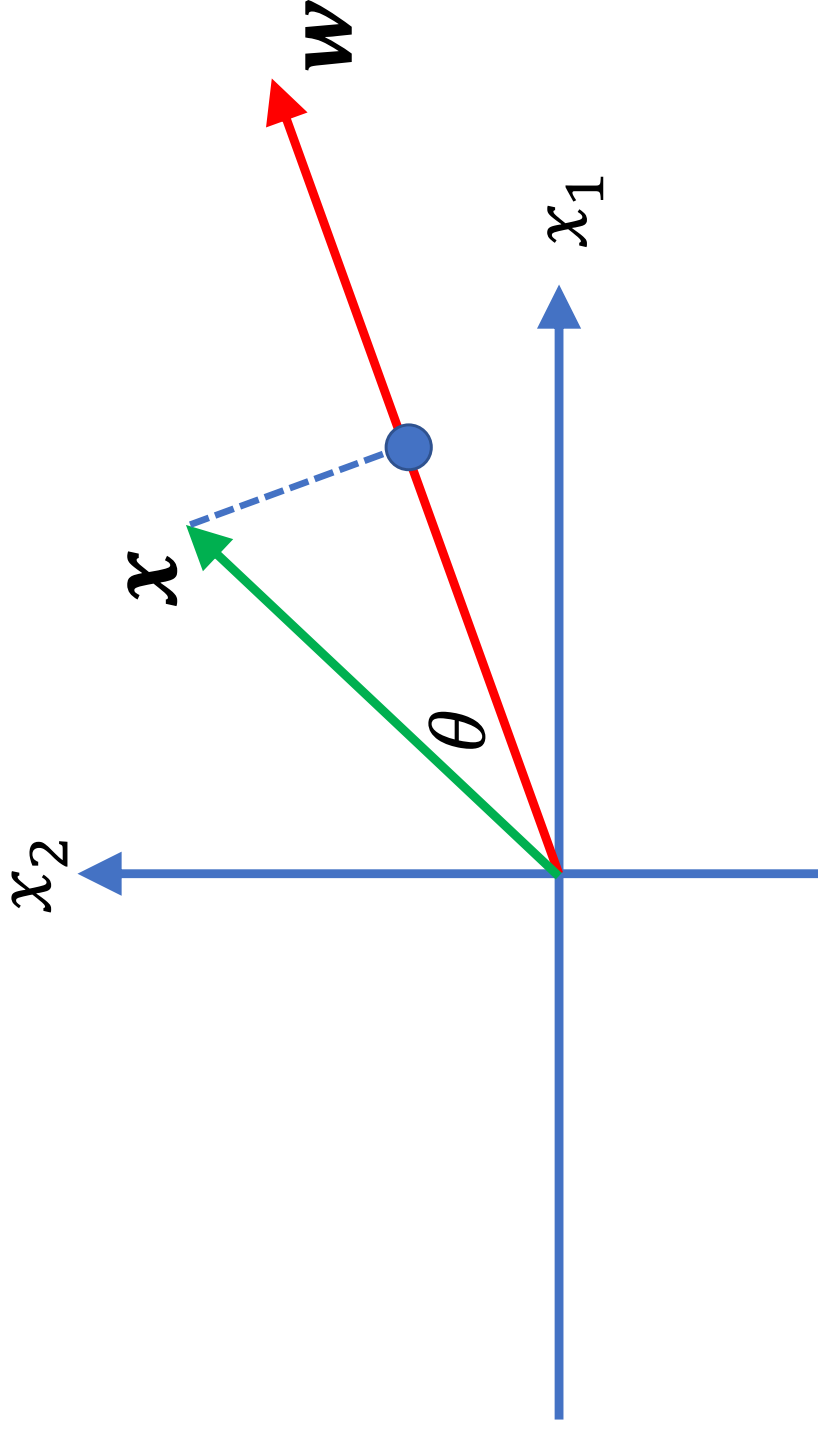
# Why not to use LS Classifier?

- Unlike LS regression, LS classification lacks a probabilistic interpretation.

- It cannot be interpreted as Maximum Likelihood of some probabilistic model on *D*.

# Fisher Discriminant Analysis (FDA)
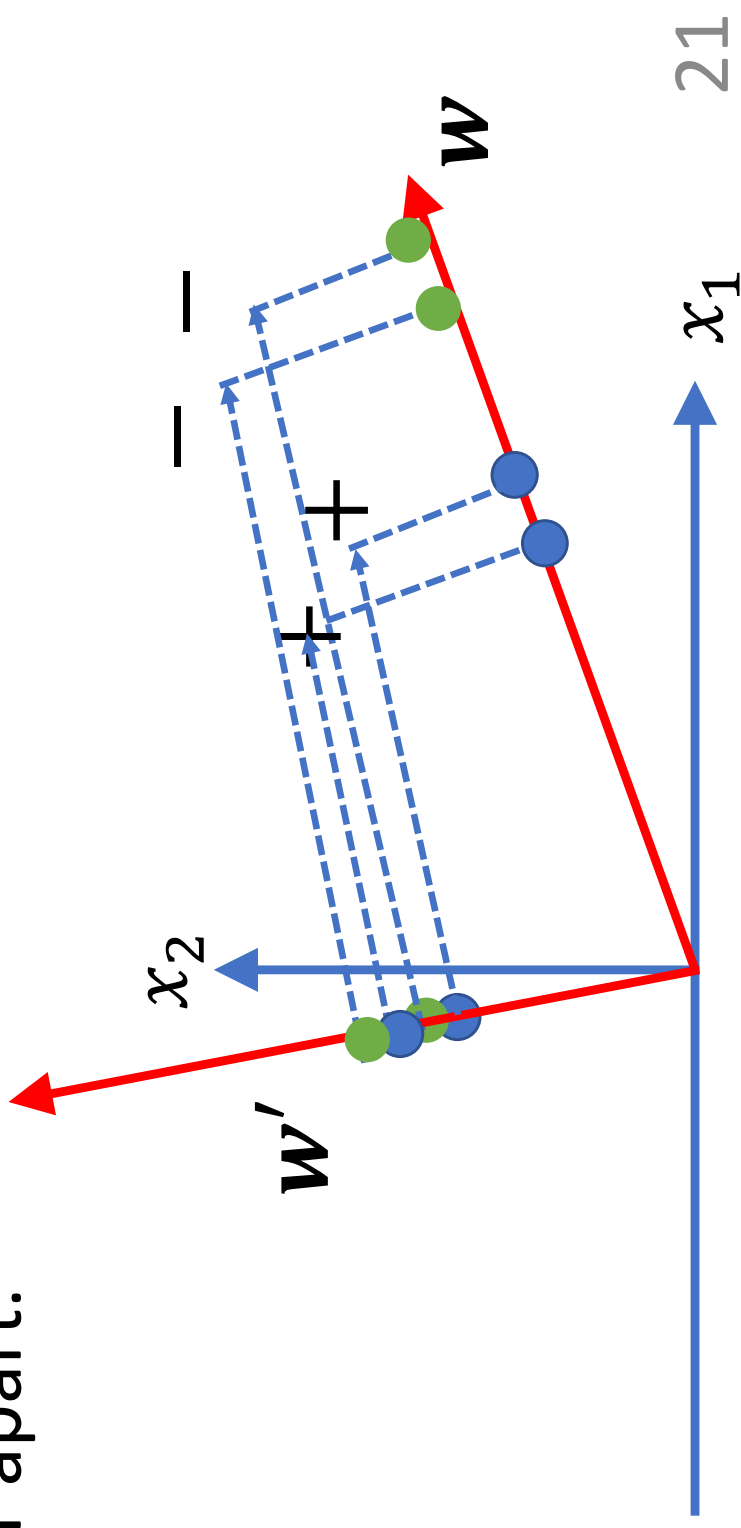
# Embedding by Inner Product

- The inner product $\langle w, x \rangle$ "embeds" $x$, onto a one-dimensional line along $w$ direction.
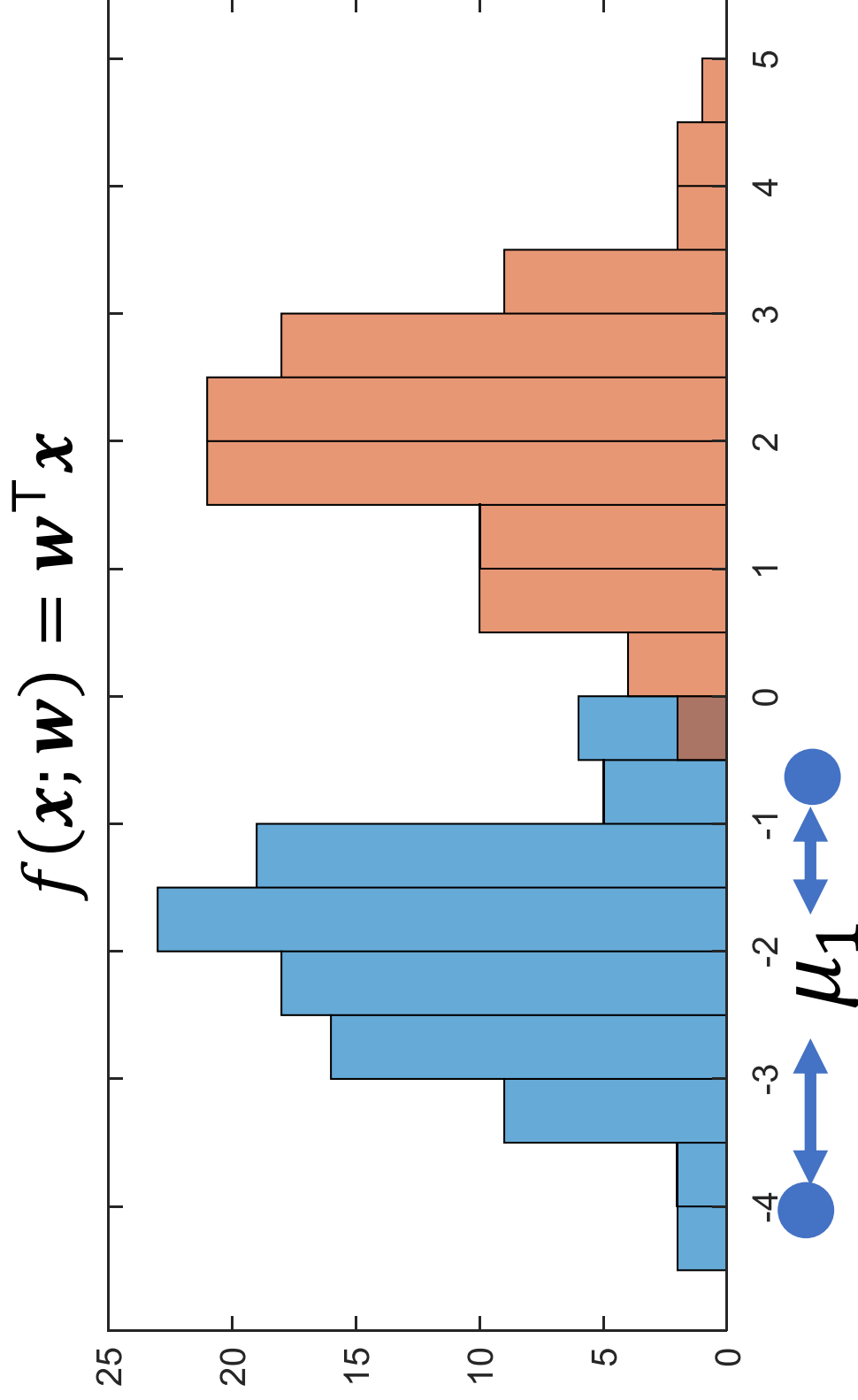
# Embedding by Inner Product

- What would be a good embedding?

- Clearly, we prefer $w$ to $w'$, as the embedding is **more separated** between $+$ and $-$ .

- We want points within the class close, but points between two classes far apart.

# Within Class and Between Class Scatterness



$f(x; w) = w^\top x$

$\mu_1 \quad \mu_1 \quad \mu$

$\mu_1 \quad \mu \quad \mu_2$

# Within-class Scatterness

- Embedding is $w^\top x$.

- Embedded center for class $k$:
  - $\hat{\mu}_k = \frac{1}{n_k} \sum_{i, y_i = k} w^\top x_i$

- Within class scatterness of class $k$:
  - $s_{w,k} = \sum_{i, y_i = k} \left( w^\top x_i - \hat{\mu}_k \right)^2$
  - Sum over points in **individual** classes.

# Between-class Scatterness

- Embedded dataset center:

  - $\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{w}^{\top} \boldsymbol{x}_i$

- Between-class scatterness
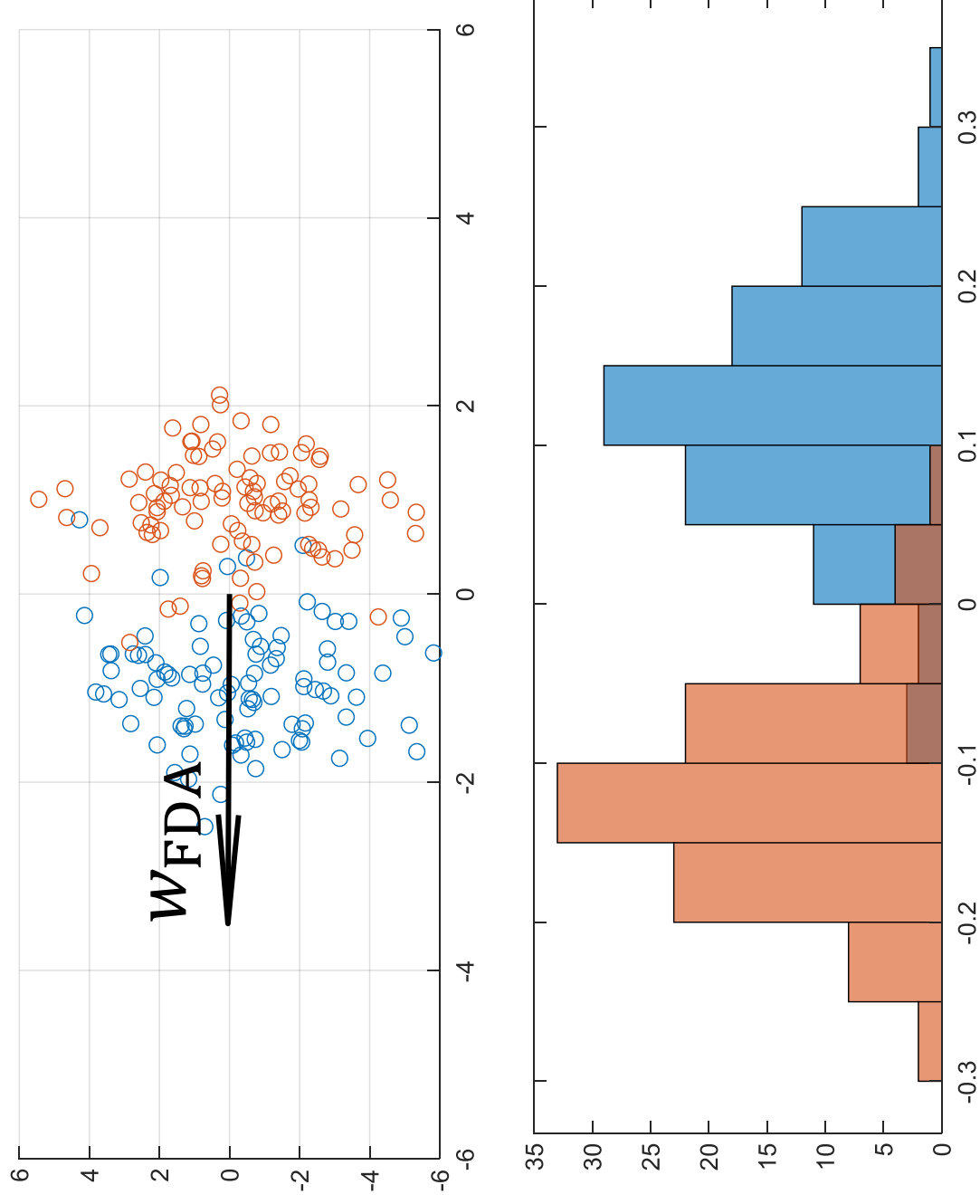
  - $s_{\mathrm{b},k} = n_k (\hat{\mu}_k - \hat{\mu})^2$

    - $n_k$ is needed to make $s_{\mathrm{b},k}$ at the same scale with $s_{\mathrm{w},k}$.

# Fisher Discriminant Analysis

- **Maximizing** between class scatterness $\forall_k$.

- **Minimize** within class scatterness $\forall_k$.

- $\max\limits_{w} \boxed{\sum_k s_{b,k}} \Big/ \boxed{\sum_k s_{w,k}}$

- If $K = 2$, this has a simple solution that

- $w := S_w^{-1}(\mu_+ - \mu_-), S_w := \sum_{k=1}^{K} S_k$

- $S_k$ is **sample covariance** of class $k$ times $n_k$.

- Read PRML 4.14 for its derivation

# Example of FDA

# Fisher Discriminant Analysis

- However, FDA does not learn a decision function $f$.

- $f(x; w_{\text{FDA}}) = \langle w_{\text{FDA}}, x \rangle$ obtained by FDA cannot be directly used for making a prediction:

- In general, $f(x; w_{\text{FDA}}) > 0$ does not mean $x$ is predicted as positive or negative data point: FDA does not care about classification accuracy, a.k.a., minimizing FP or FN.

# Probabilistic Generative Classifiers

# Probabilistic Classification

- How to put classification problem under a prob. framework?

- **Minimize Expected Loss:**

$$\hat{y} := \operatorname{argmin}_{y_0} \mathbb{E}_{p(y|x)} \left[ L(y, y_0) \,|\, x \right]$$

- We need: $p(y|x)$, $y \in \{1, \ldots, K\}$

- Discriminative: Infer $p(y|x)$ directly.
- **Generative:** Infer $p(y|x) \propto p(x|y)p(y)$, infer $p(x|y)$!

# Continuous Input Variable

- To infer $p(x|y)$, we need a model.

- If $x$ is continuous, MVN is a natural choice for $p(x|y)$.

- **Model** $p(x|y=k; w) := N_x(\mu_k, \Sigma_k)$

- <span style="color:red">Assuming IID, and all classes have shared covariance $\Sigma$</span>

- **Write down the likelihood** over $D$:

- $p(D|w) = \prod_{i \in D} p(x_i, y_i | w) = \prod_{i \in D} p(x_i | y_i; w) p(y_i)$

$$= \color{red}{\prod_{i \in D} N_{x_i}(\mu_{y_i}; \Sigma) \, p(y_i)}$$

# Continuous Input Variable

- $\widehat{\boldsymbol{\mu}}_{1\ldots K}, \widehat{\boldsymbol{\Sigma}} := \arg\max_{\boldsymbol{\mu}_{1\ldots k}, \boldsymbol{\Sigma}} \sum_{i\in D} \log[N_{x_i}(\boldsymbol{\mu}_{y_i}; \boldsymbol{\Sigma})p(y_i)]$

1. Plug in estimates for $p(y_i = k)$, which is $\frac{n_k}{n}$.

2. Now work out the MLE for $\widehat{\boldsymbol{\mu}}_k := \frac{1}{n_k}\sum_{i\in D, y_i=k} x_i$

3. Plug in $\widehat{\boldsymbol{\mu}}_k$ to work out

$$\widehat{\boldsymbol{\Sigma}} := \sum_{k=1\ldots K}\frac{n_k}{n}\boxed{\frac{1}{n_k}\sum_{i\in D,y_i=k}(\boldsymbol{x}_i - \widehat{\boldsymbol{\mu}}_k)(\boldsymbol{x}_i - \widehat{\boldsymbol{\mu}}_k)^\top}$$

MLE of covariance of individual classes!

# Linear Decision Boundary

- **Prediction:** $\hat{y} := \ \text{argmax}_y \, p(y|\boldsymbol{x}; \widehat{\boldsymbol{w}}) \propto p(\boldsymbol{x}|y; \widehat{\boldsymbol{w}}) p(y)$

- **Prove:** when using shared covariance matrix MVN model, the decision boundary is piecewise-linear.

- The decision boundary is
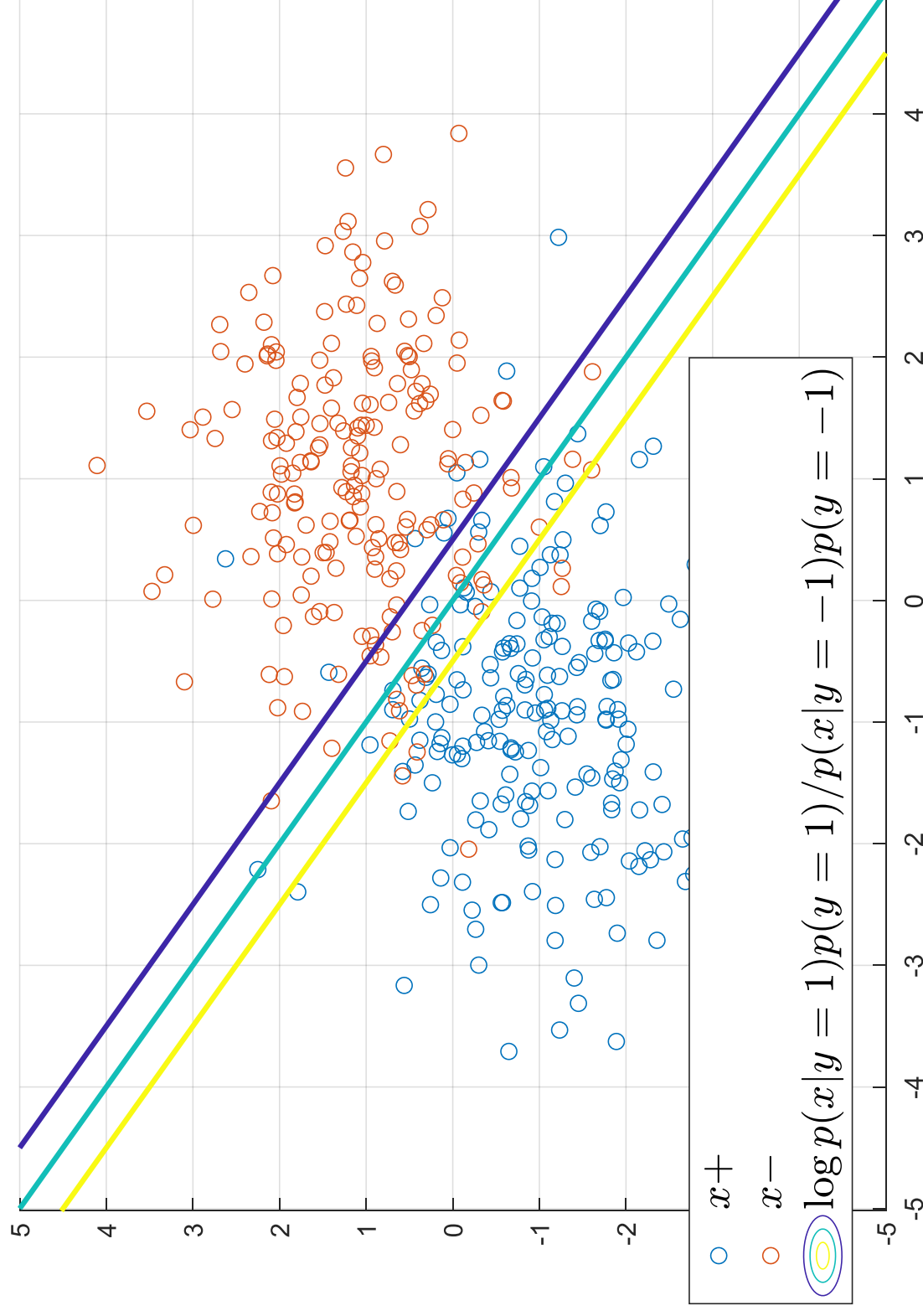$$\{\boldsymbol{x}| p(y = k|\boldsymbol{x}; \widehat{\boldsymbol{w}}) = p(y = k'|\boldsymbol{x}; \widehat{\boldsymbol{w}})\}$$
$$\forall k \neq k'$$

Which is the same as the set
$$\left\{\boldsymbol{x} \left| \frac{p(\boldsymbol{x}|y = k; \widehat{\boldsymbol{w}}) p(y = k)}{p(\boldsymbol{x}|y = k'; \widehat{\boldsymbol{w}}) p(y = k')} = 1 \right. \right\}$$
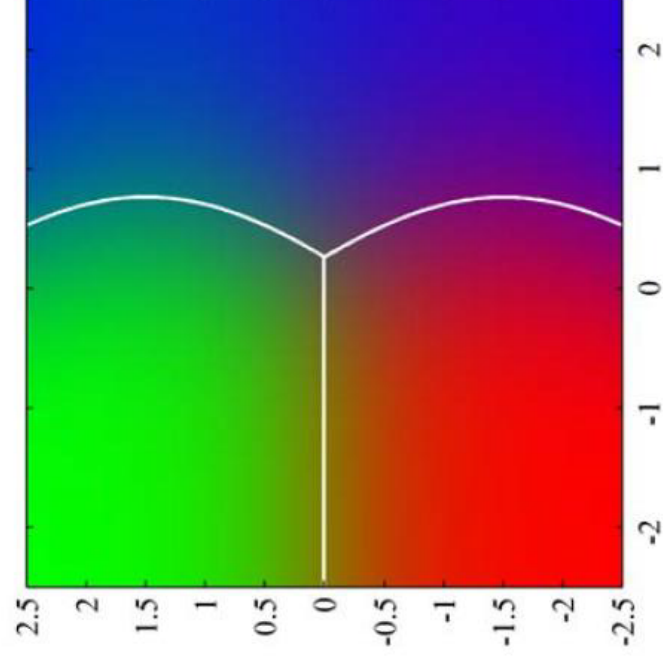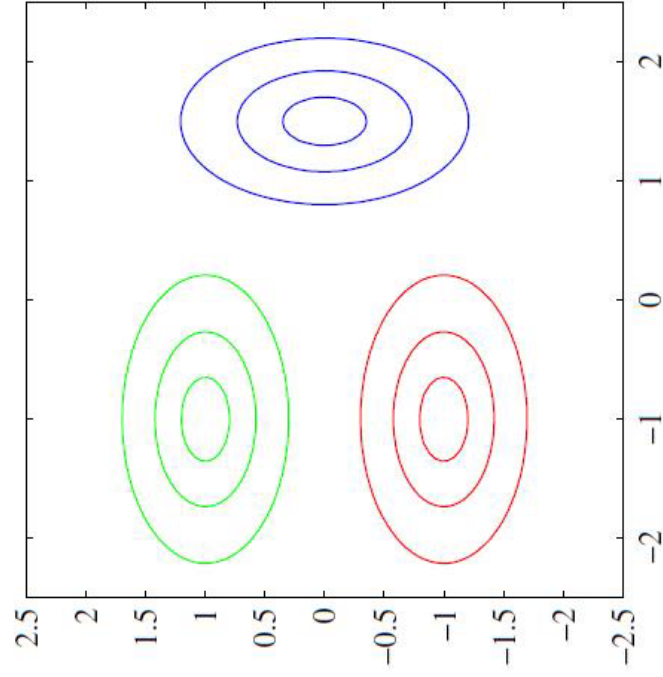$$\forall k \neq k'$$

**Hint:** take log on both sides of the equality.

# Linear Decision Boundary

# Continuous Input Variable

- You can also assume for each class $k$, there are different covariance matrices $\mathbf{\Sigma}_k$.

- The MLE reduces to estimating individual $\boldsymbol{\mu}_k$ and $\mathbf{\Sigma}_k$.

- **The decision boundary is no longer linear.**

# Discrete Input Variable $x$

- In many classification tasks, we are dealing with discrete variables as $x$. For example, in a spam filter,

- $x := \left[ x^{(1)}, \dots, x^{(d)} \right]^{\top}$ are frequencies of words in a document. This is called "bag of words" representation.

- $y \in \{spam, ham\}$.

- For example, the document "to be or not to be"

- $x := [to = 2, be = 2, or = 1, not = 1, question = 0]^{\top}$.

- $x^{(i)} \in N_0$

# Naïve Bayes

- Assume $x^{(1)} \dots x^{(d)}$ follows multinomial distribution

- $p(x = x_0|y) \propto \prod_{i=1\dots d} \beta(i|y)^{x_0^{(i)}}$ up to constant does not depend on $y$.

- $\beta(i|y = k)$ is the probability of word $i$ occurs in class $k$.

- It is easy to estimate:

$$\beta(i|y = k) \approx \frac{\sum_{j \in D, y_j = k} x_j^{(i)}}{\sum_{j \in D, y_j = k} \sum_{i=1}^{d} x_j^{(i)}}$$

- $\beta(\text{to}|y = \text{spam})$ is occurrences of the word "to" in "spam" emails divided by total number of words in "spam" emails in our training dataset.

# Naïve Bayes

- Prediction: $\hat{y} := \operatorname{argmax}_y p(\boldsymbol{x} = \boldsymbol{x}_0|y)p(y)$

- $p(y = k): \dfrac{n_k}{n}$

- $p(\boldsymbol{x} = \boldsymbol{x}_0|y) \propto \prod_{i=1\ldots d} \beta(i|y)^{x_0^{(i)}}$
  - $\beta(i|y)$ has been obtained by previous counting.

- $p(\boldsymbol{x} = \text{"}\boldsymbol{to\ be\ or\ not\ to\ be}\text{"}|y = \text{spam}) \propto$
  $\beta(to|\text{spam})^2 \beta(be|\text{spam})^2 \beta(or|\text{spam}) \beta(not|\text{spam})$

# Conclusion

- We have studied classification problem:

- Geometry of decision function

- Least square classifier

- Fisher discriminant analysis

  - Within and between scatterness

- Generative Classifiers:

  - MVN for continuous input variable

  - Naïve Bayes for discrete input variable

# Homework

- Prove the statement on page 33.

- (1) Derive the maximum likelihood estimation for parameters in multinomial distribution. (2) Explain the Naïve Bayes classifier using a Maximum Likelihood framework.

# Computing Lab

- Implement a version of Perception classifier: "Simplitron"

- Demo.