

SUPPORT VECTOR MACHINES AND KERNELIZATION

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

AN OVERVIEW

Outline of optimization steps

1. Formulate constrained form of problem
2. Convert to (**primal**) Lagrangian form
3. Take all relevant (sub)-derivatives and set to zero
4. Substitute these conditions back into **primal** Lagrangian
→ **dual** Lagrangian

(This forms a lower bound on the solution to the constrained primal form of objective)

5. Form Karush-Kuhn-Tucker (KKT) conditions for inequality constraints
6. Examine the conditions for **strong duality** which implies that the primal and dual forms have the same solution

(The usual method is via Slater's condition, which says strong duality holds if it is a convex program with non-empty constraint region)

KERNEL METHODS

INTUITION: Many methods have linear decision boundaries

We know that sometimes this isn't sufficient to represent data

EXAMPLE: Sometimes we need to include a polynomial effect or a log transform in multiple regression

Sometimes, a **linear** boundary, but in a different space makes all the difference..

OPTIMAL SEPARATING HYPERPLANE

REMINDER: The Wolfe dual, which gets maximized over α , produces the **optimal separating hyperplane**

$$\text{Wolf dual} = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k Y_i Y_k \mathbf{x}_i^\top \mathbf{x}_k$$

(this is all subject to $\alpha_i \geq 0$)

A similar result holds after the introduction of slack variables
(e.g. **support vector classifiers**)

IMPORTANT: The features only enter via

$$\mathbf{x}^\top \mathbf{x}' = \langle \mathbf{x}, \mathbf{x}' \rangle$$

(KERNEL) RIDGE REGRESSION

REMINDER: Suppose we want to predict at X , then

$$\hat{f}(X) = X^\top \hat{\beta}_{\text{ridge}}(\lambda) = X^\top \mathbb{X}^\top (\mathbb{X} \mathbb{X}^\top + \lambda I)^{-1} Y$$

Also,

$$\mathbb{X} \mathbb{X}^\top = \begin{bmatrix} \langle X_1, X_1 \rangle & \langle X_1, X_2 \rangle & \cdots & \langle X_1, X_n \rangle \\ & \vdots & & \\ \langle X_n, X_1 \rangle & \langle X_n, X_2 \rangle & \cdots & \langle X_n, X_n \rangle \end{bmatrix}$$

and

$$X^\top \mathbb{X}^\top = [\langle X, X_1 \rangle, \langle X, X_2 \rangle, \dots, \langle X, X_n \rangle]$$

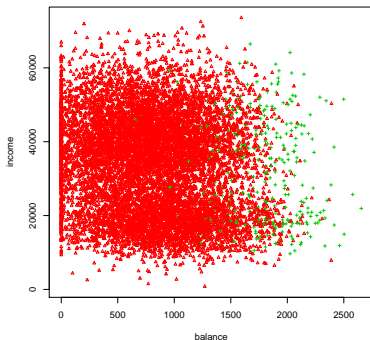
Again, we have the features enter only as

$$\langle X, X' \rangle = X^\top X'$$

LOGISTIC REGRESSION: TRANSFORMATIONS

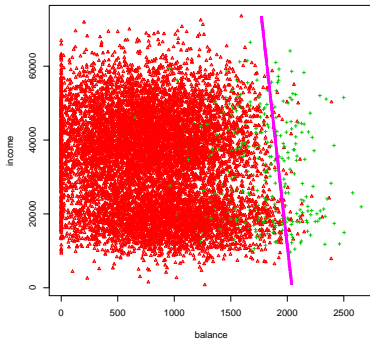
Let's look at the **default** data in “Introduction to Statistical Learning”

In particular, we will look at **default** status as a function of **balance** and **income**



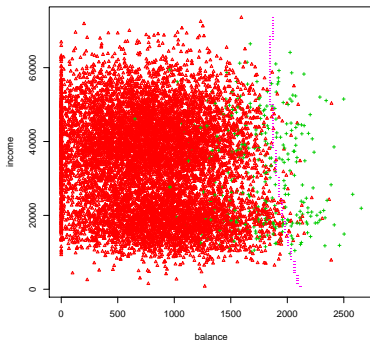
LOGISTIC REGRESSION: TRANSFORMATIONS

```
out.glm = glm(default~balance + income,family='binomial')
```



LOGISTIC REGRESSION: TRANSFORMATIONS

```
out.glm = glm(default~balance + income +  
              I(income^2),family='binomial')
```



CONCLUSION: A **Linear** rule in a transformed space can have a **nonlinear** boundary in the original features

LOGISTIC REGRESSION: TRANSFORMATIONS

REMINDER: The logistic model: untransformed

$$\begin{aligned}\text{logit}(\mathbb{P}(Y = 1|X)) &= \beta_0 + \beta^\top X \\ &= \beta_0 + \beta_1 \text{balance} + \beta_2 \text{income}\end{aligned}$$

The decision boundary is the hyperplane $\{X : \beta_0 + \beta^\top X = 0\}$

This is **linear** in the feature space

LOGISTIC REGRESSION: TRANSFORMATIONS

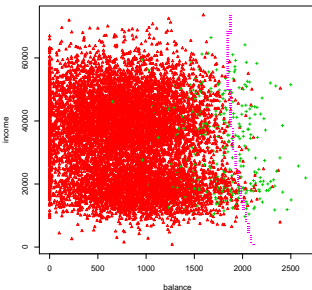
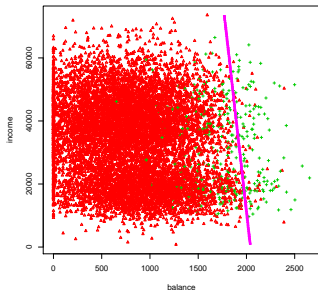
Adding the polynomial transformation $\Phi(X) = (X_1, X_2, X_2^2)$:

$$\text{logit}(\mathbb{P}(Y = 1|X)) = \beta_0 + \beta^\top \Phi(X)$$

$$= \beta_0 + \beta_1 \text{balance} + \beta_2 \text{income} + \beta_3 \text{income}^2$$

Decision boundary is still a hyperplane $\{X : \beta_0 + \beta^\top \Phi(X) = 0\}$

This is **nonlinear** in the feature space!



LOGISTIC REGRESSION: TRANSFORMATIONS

Of course, as we include more transformations,

- We need to choose the transformations **manually**
- **Computations** can become difficult if we aren't careful
(**EXAMPLE:** Solving the least squares problem takes something like np^2 computations)
- We need to **regularize** to prevent overfitting

Can we form them in an automated fashion?

Kernel Methods

THREE RELATED METHODS

The following are seemingly disparate methods

- **SMOOTHNESS PENALIZATION:** Regularizing a loss function with a penalty on smoothness
(Example: Smoothing splines)
- **FEATURE CREATION:** Imposing a **feature mapping** $\Phi : \mathbb{R}^p \rightarrow \mathcal{A}$ thus creating new features e.g. via polynomials or interactions
(Example: Regression splines or polynomial regression)
- **GAUSSIAN PROCESSES:** Modeling the regression function as a Gaussian process with a given mean and covariance
(Example: Gaussian process regression)

It turns out these concepts are all the same and each forms a **reproducing kernel Hilbert space** (RKHS)

(Many of these ideas are in Wahba (1990). It was introduced to the ML community in Vapnik et al. (1996) and summarized in a nice review paper in Hofmann et al. (2008))

KERNEL METHODS

Suppose $k : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$ is a **positive definite** kernel

(This means $\int \int k(x, y) f(x) f(y) dx dy > 0$)

To be concrete, think of $\mathcal{A} = \mathbb{R}^p$

(However, any set of objects will do as long as an **inner product** can be defined)

Let's consider the space of functions generated by the completion of

$$\mathcal{H}_k = \{k(\cdot, y) : y \in \mathbb{R}^p\}$$

(This, loosely speaking, is all functions of the form $f(x) = \sum_{j=1}^J \alpha_j k(x, y_j)$)

NONNEGATIVE DEFINITE MATRICES

Let $A \in \mathbb{R}^{p \times p}$ be a symmetric, nonnegative definite matrix:

$$z^\top A z \geq 0 \text{ for all } z \text{ and } A^\top = A$$

Then, A has an eigenvalue expansion

$$A = UDU^\top = \sum_{j=1}^p d_j u_j u_j^\top$$

where $d_j \geq 0$

OBSERVATION: Each such A , generates a new inner product

$$\langle z, z' \rangle = z^\top z' = z^\top \underbrace{I}_{\text{Identity}} z'$$

$$\langle z, z' \rangle_A = z^\top A z'$$

(If we enforce A to be positive definite, then $\langle z, z \rangle_A = \|z\|_A^2$ is a norm)

NONNEGATIVE DEFINITE MATRICES

Suppose A_{ij}^j is the (i, j) entry in A , and A_i is the i^{th} row

$$Az = \begin{bmatrix} A_1^\top \\ \vdots \\ A_p^\top \end{bmatrix} z = \begin{bmatrix} A_1^\top z \\ \vdots \\ A_p^\top z \end{bmatrix}$$

NOTE: Multiplication by A is really taking **inner products** with its rows.

Hence, A_i is called the (multiplication) **kernel** of matrix A

KERNEL METHODS

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a **symmetric, nonnegative definite** kernel

Write the eigenvalue expansion of k as

$$k(X, X') = \sum_{j=1}^{\infty} \theta_j \phi_j(X) \phi_j(X')$$

- $\theta_j \geq 0$ (**nonnegative definite**)
- $\|(\theta_j)_{j=1}^{\infty}\|_2^2 = \sum_{j=1}^{\infty} \theta_j^2 < \infty$
- The ϕ_j are orthogonal **eigenfunctions**: $\int \phi_j \phi_{j'} = \delta_{j,j'}$

(This is called **Mercer's theorem**, and such a k is called a **Mercer** kernel)

We can write any $f \in \mathcal{H}_k$ with two constraints

- $f(X) = \sum_{j=1}^{\infty} f_j \phi_j(X)$
- $\langle f, f \rangle_{\mathcal{H}_k} = \|f\|_{\mathcal{H}_k}^2 = \sum_{j=1}^{\infty} f_j^2 / \theta_j < \infty$

KERNEL METHODS VIA REGULARIZATION

After specifying a kernel function¹ k , we can define an estimator via

$$\min_{f \in \mathcal{H}_k} \hat{\mathbb{P}} \ell_f + \lambda \|f\|_{\mathcal{H}_k}^2$$

This is a (potentially) infinite dimensional optimization problem

(hard, especially with a computer)

It can be shown that the solution has the form

$$f(X) = \sum_{i=1}^n \beta_i k(X, X_i)$$

(This is known as the **representer theorem**)

¹Or crucially and equivalently a set of eigenfunctions and eigenvalues

KERNEL METHODS VIA REGULARIZATION

$$f(X) = \sum_{i=1}^n \beta_i k(X, X_i)$$

The terms $k(X, X_i)$ are the **representers**, as

$$\langle k(\cdot, X_i), f \rangle_{\mathcal{H}_k} = f(X_i)$$

and \mathcal{H}_k is called a **reproducing kernel Hilbert space** (RKHS) as

$$\langle k(\cdot, X_i), k(\cdot, X_{i'}) \rangle_{\mathcal{H}_k} = k(X_i, X_{i'})$$

KERNEL METHODS VIA REGULARIZATION

Due to these properties, we can write the optimization problem as

$$\min_{\beta} \hat{\mathbb{P}} \ell_{\mathbb{K}} \beta + \lambda \beta^{\top} \mathbb{K} \beta$$

where $\mathbb{K} = [k(X_i, X_{i'})]$

This provides a prescription for forming an incredibly rich suite of estimators:

Choose a

- kernel k
- loss function ℓ

and then minimize

KERNEL METHODS VIA REGULARIZATION: EXAMPLE

Suppose that $\ell_{\mathbb{K}\beta}(Z) = (Y - \mathbb{K}\beta)^2$

Then:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \hat{\mathbb{P}}\ell_{\mathbb{K}\beta} + \lambda\beta^{\top}\mathbb{K}\beta = (\mathbb{K} + \lambda I)^{-1}Y$$

and

$$\hat{f} = \mathbb{K}\hat{\beta} = \mathbb{K}(\mathbb{K} + \lambda I)^{-1}Y = (\lambda\mathbb{K}^{-1} + I)^{-1}Y$$

are the **fitted values**

(This should be compared with the notes on ridge regression)

KERNEL: EXAMPLE

Back to polynomial terms/interactions:

Form

$$k_d(X, X') = (X^\top X' + 1)^d$$

k_d has $M = \binom{p+d}{d}$ eigenfunctions

These **span** the space of polynomials in \mathbb{R}^p with degree d

KERNEL: EXAMPLE

EXAMPLE: Let $d = p = 2 \Rightarrow M = 6$ and

$$\begin{aligned}k(u, v) &= 1 + 2u_1v_1 + 2u_2v_2 + u_1^2v_1^2 + u_2^2v_2^2 + 2u_1u_2v_1v_2 \\&= \sum_{k=1}^M \Phi_k(u)\Phi_k(v) \\&= \Phi(u)^\top \Phi(v) \\&= \langle \Phi(u), \Phi(v) \rangle\end{aligned}$$

where

$$\Phi(v)^\top = (1, \sqrt{2}v_1, \sqrt{2}v_2, v_1^2, v_2^2, \sqrt{2}v_1v_2)$$

IMPORTANT: These equalities are **everything** that makes kernelization work!

KERNEL: CONCLUSION

Let's recap:

$$\begin{aligned}k(u, v) &= 1 + 2u_1v_1 + 2u_2v_2 + u_1^2v_1^2 + u_2^2v_2^2 + 2u_1u_2v_1v_2 \\ &= \langle \Phi(u), \Phi(v) \rangle\end{aligned}$$

- Some methods only involve features via inner products $X^\top X' = \langle X, X' \rangle$
(We've explicitly seen two: ridge regression and support vector classifiers)
- If we make transformations of X to $\Phi(X)$, the procedure depends on $\Phi(X)^\top \Phi(X') = \langle \Phi(X), \Phi(X') \rangle$
- **CRUCIAL:** We can compute this inner product via the kernel:

$$k(X, X') = \langle \Phi(X), \Phi(X') \rangle$$

KERNEL: CONCLUSION

Instead of creating a very high dimensional object via transformations, choose a kernel k

Now, the only thing left to do is form the **outer product** of kernel evaluations

$$\mathbb{K} = [k(X_i, X_{i'})]_{1 \leq i, i' \leq n}$$

```
x = c(1,2,3)# n = 3
k = function(x,y){ return(x + y + x*y)}
> outer(x,x,k)
      [,1] [,2] [,3]
[1,]    3    5    7
[2,]    5    8   11
[3,]    7   11   15
```

(Kernel) SVMs

KERNEL SVM

RECALL:

$$\frac{1}{2} \|\beta\|_2^2 - \sum_{i=1}^n \alpha_i [Y_i (X_i^\top \beta + \beta_0) - 1]$$

Derivatives with respect to β and β_0 imply:

- $\beta = \sum_{i=1}^n \alpha_i Y_i X_i$
- $0 = \sum_{i=1}^n \alpha_i Y_i$

Write the solution function

$$h(X) = \beta_0 + \beta^\top X = \beta_0 + \sum_{i=1}^n \alpha_i Y_i X_i^\top X$$

Kernelize the support vector classifier \Rightarrow support vector machine (SVM):

$$h(X) = \beta_0 + \sum_{i=1}^n \alpha_i Y_i k(X_i, X)$$

GENERAL KERNEL MACHINES

After specifying a kernel function, it can be shown that many procedures have a solution of the form

$$f(X) = \sum_{i=1}^n \gamma_i k(X, X_i)$$

For some $\gamma_1, \dots, \gamma_n$

Also, this is equivalent to performing the method in the space given by the **eigenfunctions** of k

$$k(u, v) = \sum_{j=1}^{\infty} \theta_j \phi_j(u) \phi_j(v)$$

Also, (the) **feature map** is

$$\Phi = [\phi_1, \dots, \phi_p, \dots]$$

KERNEL SVMs

Hence (and luckily) specifying Φ itself unnecessary,

(Luckily, as many kernels have difficult to compute eigenfunctions)

We need only define the **kernel** that is symmetric, positive definite

Some common choices for SVMs:

- **POLYNOMIAL:** $k(x, y) = (1 + x^\top y)^d$
- **RADIAL BASIS:** $k(x, y) = e^{-\tau \|x - y\|_b^b}$

(For example, $b = 2$ and $\tau = 1/(2\sigma^2)$ is (proportional to) the Gaussian density)

KERNEL SVMs: SUMMARY

Reminder: the solution form for SVM is

$$\beta = \sum_{i=1}^n \alpha_i Y_i X_i$$

Kernelized, this is

$$\beta = \sum_{i=1}^n \alpha_i Y_i \Phi(X_i)$$

Therefore, the induced hyperplane is:

$$\begin{aligned} h(X) &= \Phi(X)^\top \beta + \beta_0 = \sum_{i=1}^n \alpha_i Y_i \langle \Phi(X), \Phi(X_i) \rangle + \beta_0 \\ &= \sum_{i=1}^n \alpha_i Y_i k(X, X_i) + \beta_0 \end{aligned}$$

The final classification is still $\hat{g}(X) = \text{sgn}(\hat{h}(X))$

SVMs via penalization

SVMs VIA PENALIZATION

NOTE: SVMs can be derived from **penalized loss** methods

The support vector classifier optimization problem:

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|_2^2 + \lambda \sum \xi_i \quad \text{subject to}$$

$$Y_i h(X_i) \geq 1 - \xi_i, \xi_i \geq 0, \text{ for each } i$$

Consider the alternative program:

$$\min_{\beta_0, \beta} \sum_{i=1}^n [1 - Y_i h(X_i)]_+ + \tau \|\beta\|_2^2$$

These optimization problems are the same!

(With the relation: $2\lambda = 1/\tau$)

SVMs VIA PENALIZATION

The **loss** part is the **hinge loss function**

$$\ell(X, Y) = [1 - Yh(X)]_+$$

The hinge loss approximates the zero-one loss function underlying classification

It has one major advantage, however: **convexity**

SURROGATE LOSSES: CONVEX RELAXATION

Looking at

$$\min_{\beta, \beta_0} \sum_{i=1}^n [1 - Y_i h(X_i)]_+ + \tau \|\beta\|_2^2$$

It is tempting to minimize (analogous to linear regression)

$$\sum_{i=1}^n \mathbf{1}(Y_i \neq \hat{g}(X_i)) + \tau \|\beta\|_2^2$$

However, this is **nonconvex** (in $u = h(X)Y$)

A common trick is to approximate the **nonconvex** objective with a convex one

(This is known as **convex relaxation** with a **surrogate loss function**)

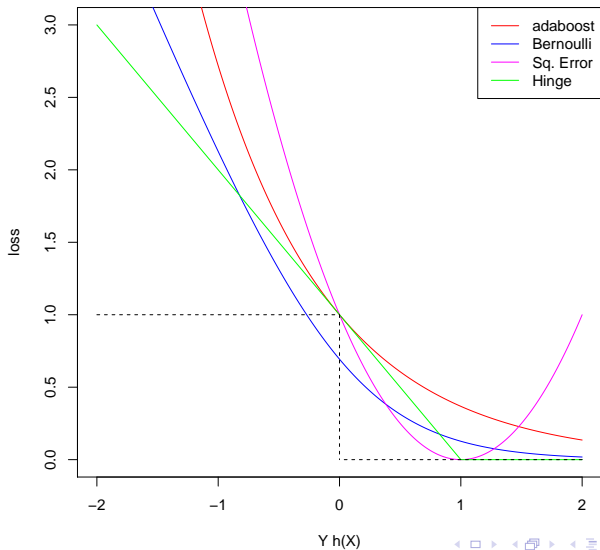
SURROGATE LOSSES

IDEA: We can use a **surrogate** loss that mimics this function while still being convex

It turns out we have already done that! (three times)

- **HINGE:** $[1 - Yh(X)]_+$
- **LOGISTIC:** $\log(1 + e^{-Yh(X)})$
- **ADABOOST:** $e^{-Yh(X)}$

COMPARING LOSS FUNCTIONS



SVMs IN PRACTICE

GENERAL FUNCTIONS: The basic SVM functions are in the C++ library **libsvm**

R PACKAGE: The **R** package **e1071** calls **libsvm**

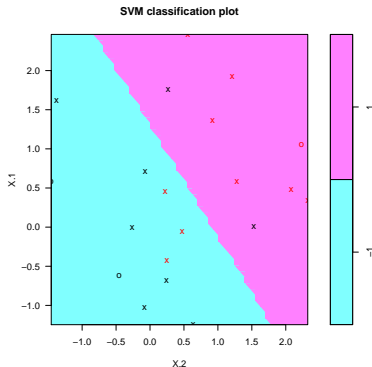
PATH ALGORITHM: **svmpath**

For a nice comparison of these approaches, see “Support vector machines in **R**”

(<http://www.jstatsoft.org/v15/i09/paper>)

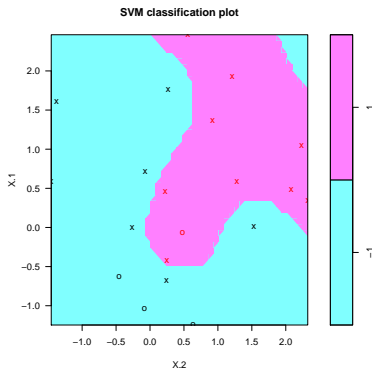
SVM EXAMPLE

```
tune.out = tune(svm,Y~.,data=dat,kernel="linear",  
               ranges=list(cost=c(0.001, 0.01, 0.1, 1,5,10,100)))
```



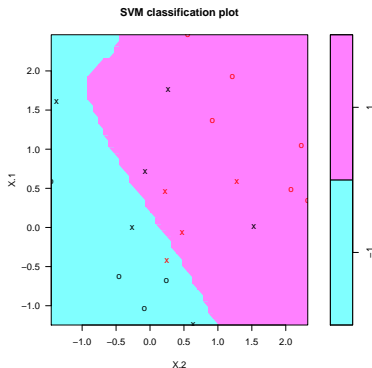
SVM EXAMPLE

```
tune.out = tune(svm,Y~.,data=dat,kernel="radial",  
  gamma=c(1,2),  
  ranges=list(cost=c(0.001, 0.01, 0.1, 1,5,10,100)))
```

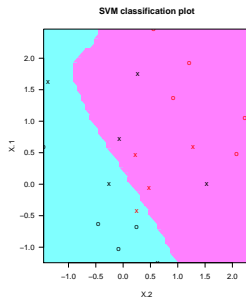
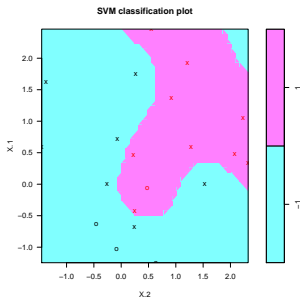
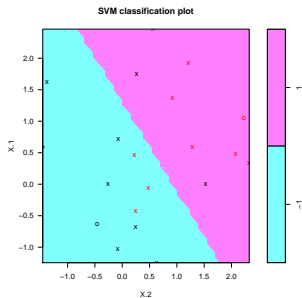


SVM EXAMPLE

```
tune.out = tune(svm,Y~.,data=dat,kernel="polynomial",  
  degree=c(3,5,10),  
  ranges=list(cost=c(0.001, 0.01, 0.1, 1,5,10,100)))
```



SVM EXAMPLE



Multiclass classification

MULTICLASS SVMs

Sometimes, it becomes necessary to do multiclass classification

There are two main approaches:

- One-versus-one
- One-versus-all

MULTICLASS SVMs: ONE-VERSUS-ONE

Here, for G possible classes, we run $G(G - 1)/2$ possible pairwise classifications

For a given test point X , we find $\hat{g}_k(X)$ for $k = 1, \dots, G(G - 1)/2$ fits

The result is a vector $\hat{G} \in \mathbb{R}^G$ with the total number of times X was assigned to each class

We report $\hat{g}(X) = \arg \max_g \hat{G}$

This approach uses all the class information, but can be slow

MULTICLASS SVMs: ONE-VERSUS-ALL

Here, we fit only G SVMs by respectively collapsing over all size $G - 1$ subsets of $\{1, \dots, G\}$

(This is compared with $G(G - 1)/2$ comparisons for one-versus-one)

Take all $\hat{h}_g(X)$ for $g = 1, \dots, G$, where class g is coded 1 and “the rest” is coded -1

Assign $\hat{g}(X) = \arg \max_g \hat{h}_g(X)$

(Note that these strategies can be applied to any classifier)