

Support Vector Machines 2; Multiclass Classification

Λ

Linear

Outline

- SVM Review
- Support vectors
- Multiclass linear classificaiton

SVM Motivation

*Tabular
Data*

- The support vector machine (SVM) is one of the best ``off-the-shelf'' classification methods.
- “ ... looking at Kaggle challenges that are not related to vision or sequential tasks, gradient boosting, random forests, or SVMs are winning most of the competitions” (Klambauer et al., NeurIPS 2017)
- “On 75 small datasets with less than 1000 data points, random forests and SVMs outperform [FNNs]. On 46 larger datasets with at least 1000 data points, SNNs show the highest performance followed by SVMs and random forests (Klambauer et al., NeurIPS 2017)
- "We evaluate 179 classifiers arising from 17 families ... We use 121 data sets ... The classifiers most likely to be the best are the random forest (RF) versions ... However, the difference is not statistically significant with the second best, the SVM with Gaussian kernel implemented in C using LibSVM“ (Delgado et al., JMLR 2014)

OSM Hyperplane

- The SVM is obtained by kernelizing the OSM hyperplane (equivalently, regularized ERM with the hinge loss).

$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \quad (\text{OSM})$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$



$$(\mathbf{w}^*, b^*, \xi^*)$$

$$x \mapsto \text{sign} \left\{ (\mathbf{w}^*)^T x + b^* \right\}$$

Poll

The number of optimization variables in the dual optimization problem (associated to the OSM hyperplane) is on the order of

- (A) d , the original input dimension
- (B) m , the dimensions of the output of the feature map Φ
- (C) n , the number of training data points
- (D) None of the above

OSM Dual Formulation

- The OSM hyperplane can be kernelized by solving the dual

$$\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq \frac{C}{n}, \quad \forall i = 1, \dots, n \end{aligned}$$

- The classifier is expressed

$$f(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b^* \right\}$$

- Taking any j with $0 < \alpha_j^* < \frac{C}{n}$, the offset is given by

$$b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

The Support Vector Machine

- Let k be an inner product/SPD kernel $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$
- The *support vector machine* is the classifier obtained by solving

$$\max_{\boldsymbol{\alpha}} -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq \frac{C}{n}, \quad \forall i = 1, \dots, n$$

- The classifier is expressed

$$f(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^* \right\}$$

where $b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}_j)$ for any j with $0 < \alpha_j^* < \frac{C}{n}$.

Support Vectors

- From complementary slackness,

$$H_i - \alpha_i^* \left(1 - \xi_i^* - y_i (\langle w^*, x_i \rangle + b^*) \right) = 0$$

$\underbrace{\qquad\qquad\qquad}_{g_i}$

- If x_i satisfies

$$y_i (\langle w^*, x_i \rangle + b^*) = 1 - \xi_i^*$$

we call x_i a *support vector*.

- Therefore, if x_i is *not* a support vector, then $\alpha_i^* = 0$
- Conclusion: The SVM classifier depends only on the SVs:

$$f(x) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i k(x, x_i) + b^* \right\} = \text{sign} \left\{ \sum_{\substack{i: x_i \\ \text{is a SV}}} \alpha_i^* y_i k(x, x_i) + b^* \right\}$$

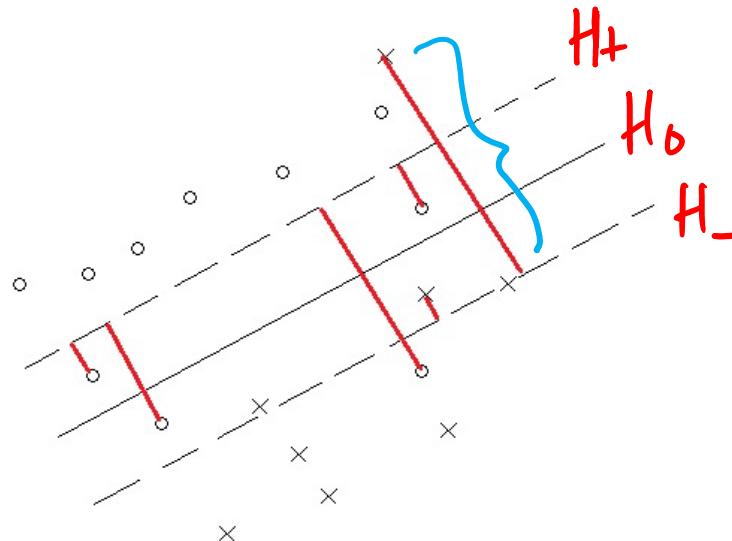
Margin Errors

- Recall the *decision boundary* is the hyperplane

$$H_0 := \{x \mid \langle w^*, x \rangle + b^* = 0\}$$

- The *margin boundaries* are hyperplanes H_+ and H_- defined by

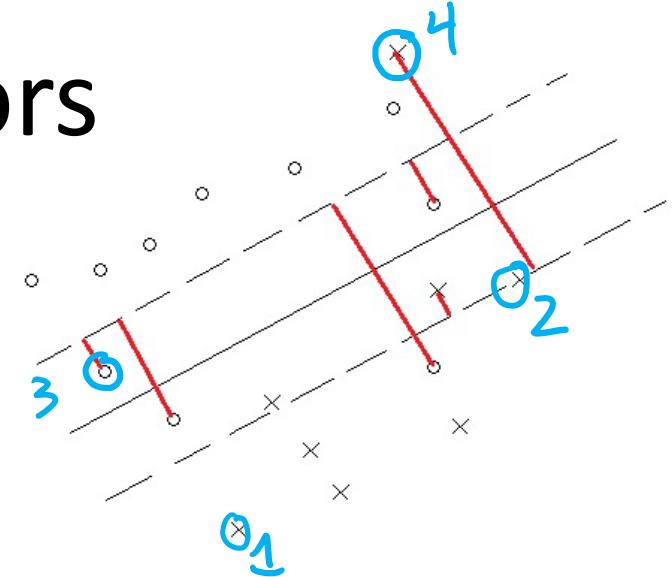
$$H_{\pm} := \{x \mid \langle w^*, x \rangle + b^* = \pm 1\}$$



x_i is SV
↑

Margin Errors

$$1 - \xi_i^* - y_i (\langle w^*, x_i \rangle + b^*) = 0$$



- Optimal slack variables satisfy

$$\xi_i^* = \max(0, 1 - y_i((w^*)^T x + b^*)).$$

- Consider four cases:

- If $y_i(\langle w^*, x_i \rangle + b^*) > 1$, then $\xi_i^* = 0$ and x_i is not a support vector.
- If $y_i(\langle w^*, x_i \rangle + b^*) = 1$, then $\xi_i^* = 0$ and x_i is a support vector.
- If $0 \leq y_i(\langle w^*, x_i \rangle + b^*) < 1$, then $\xi_i^* > 0$ and x_i is a support vector. "within the margin"
- If $y_i(\langle w^*, x_i \rangle + b^*) < 0$, then $\xi_i^* > 1$ and x_i is a support vector.

- Case 2-4 are referred to as *margin errors* - these are the SV

SVM Margin Errors

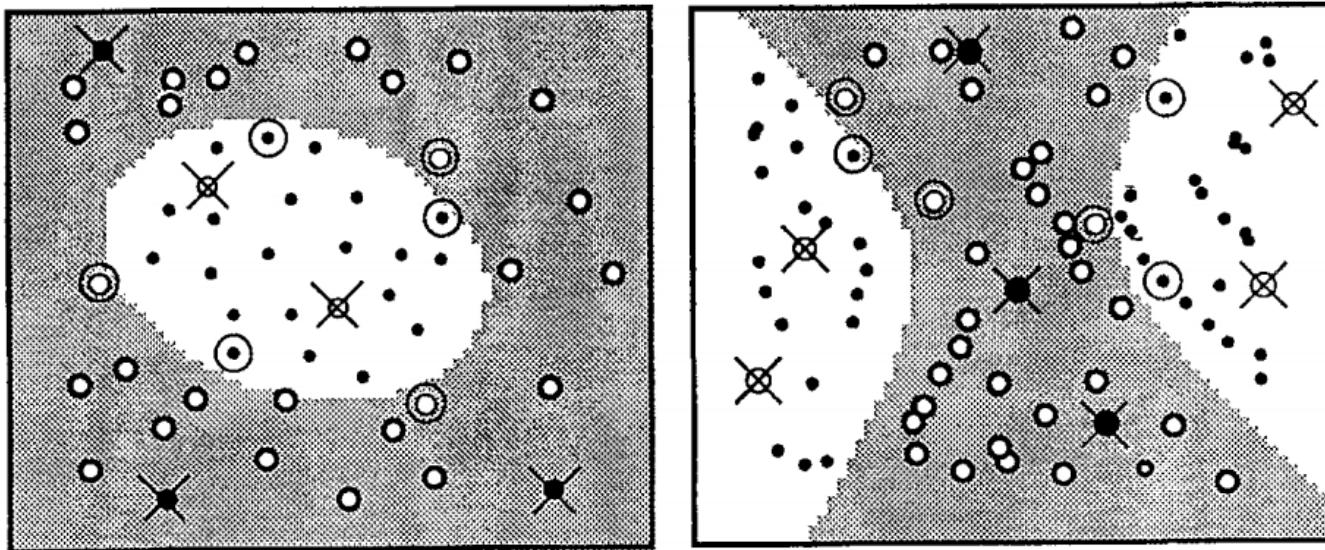


Figure from the original paper by Corinna Cortes and Vladimir Vapnik,
“Support vector networks”, *Machine Learning Journal*, 1995.

Poll

True or false: If x_i is not a support vector and it is removed from the training data set, the SVM classifier does not change.

(A) True

(B) False

Remarks

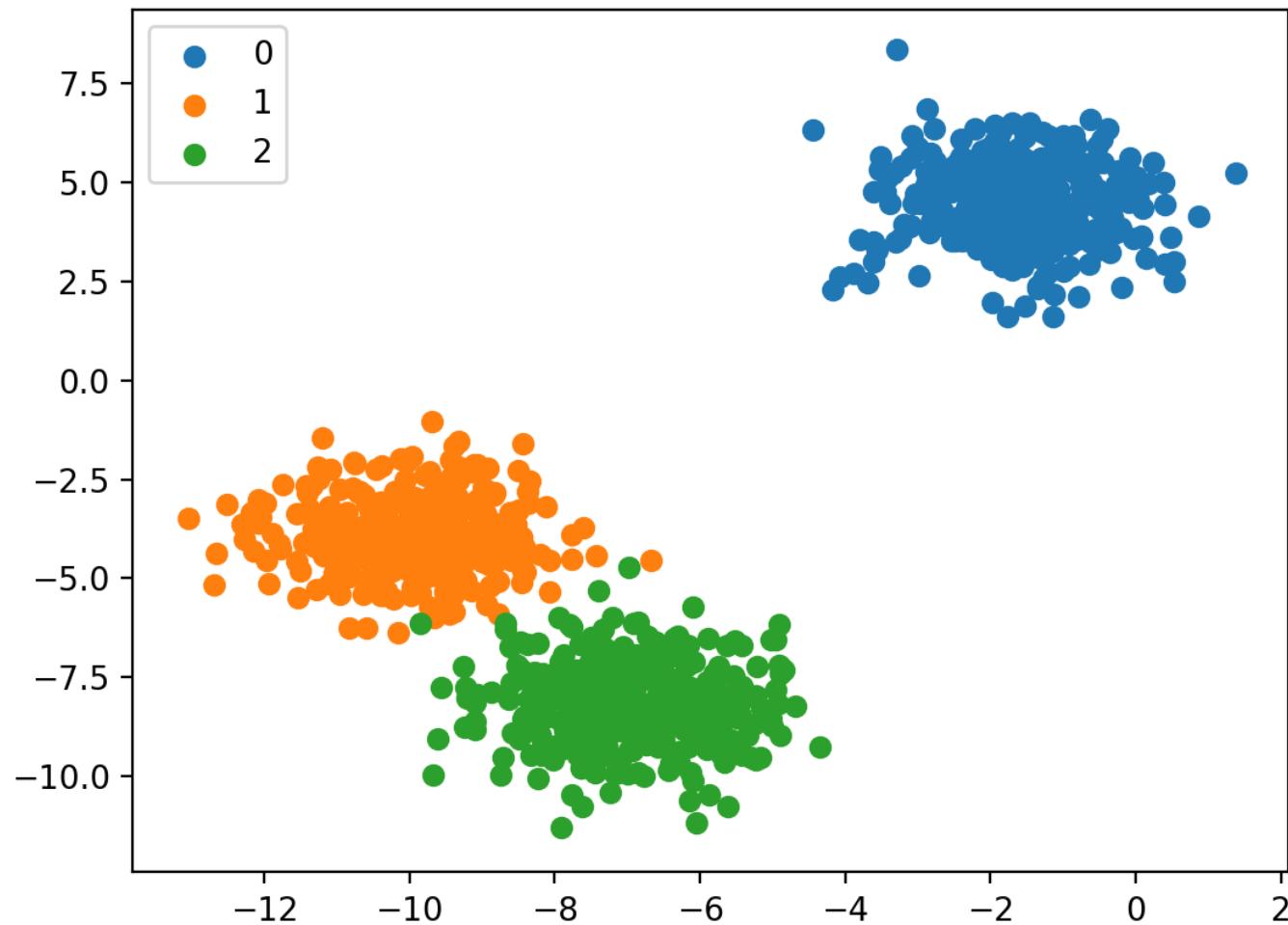
- The final classifier depends only on those training data points that are support vectors. In practice this is often only a small fraction of the original training data, so the classifier can be stored and evaluated efficiently.
- The size of the dual (i.e., the number of variables) is n , and in particular it is independent of the output dimension of the feature map Φ associated with k , which could be infinite.
- Computational complexity of solving dual is $\Omega(dn^2)$ and $O(dn^3)$.
- Modern solvers use (block) coordinate descent and can handle data sets with tens of millions of points

Remark

- Suppose $\boldsymbol{\alpha}^*$ is a dual solution, and that \mathbf{x}_i is not a support vector for a certain i . Then $\alpha_i^* = 0$.
- Now consider the modified data set obtained by removing (\mathbf{x}_i, y_i) . Let $\boldsymbol{\alpha}_{-i}^*$ be the length $n - 1$ vector obtained by removing the i -th entry from $\boldsymbol{\alpha}^*$.
- Then $\boldsymbol{\alpha}_{-i}^*$ solves the dual for the modified dataset.
- Therefore the SVM remains unchanged.

Multiclass Classification

Lihear



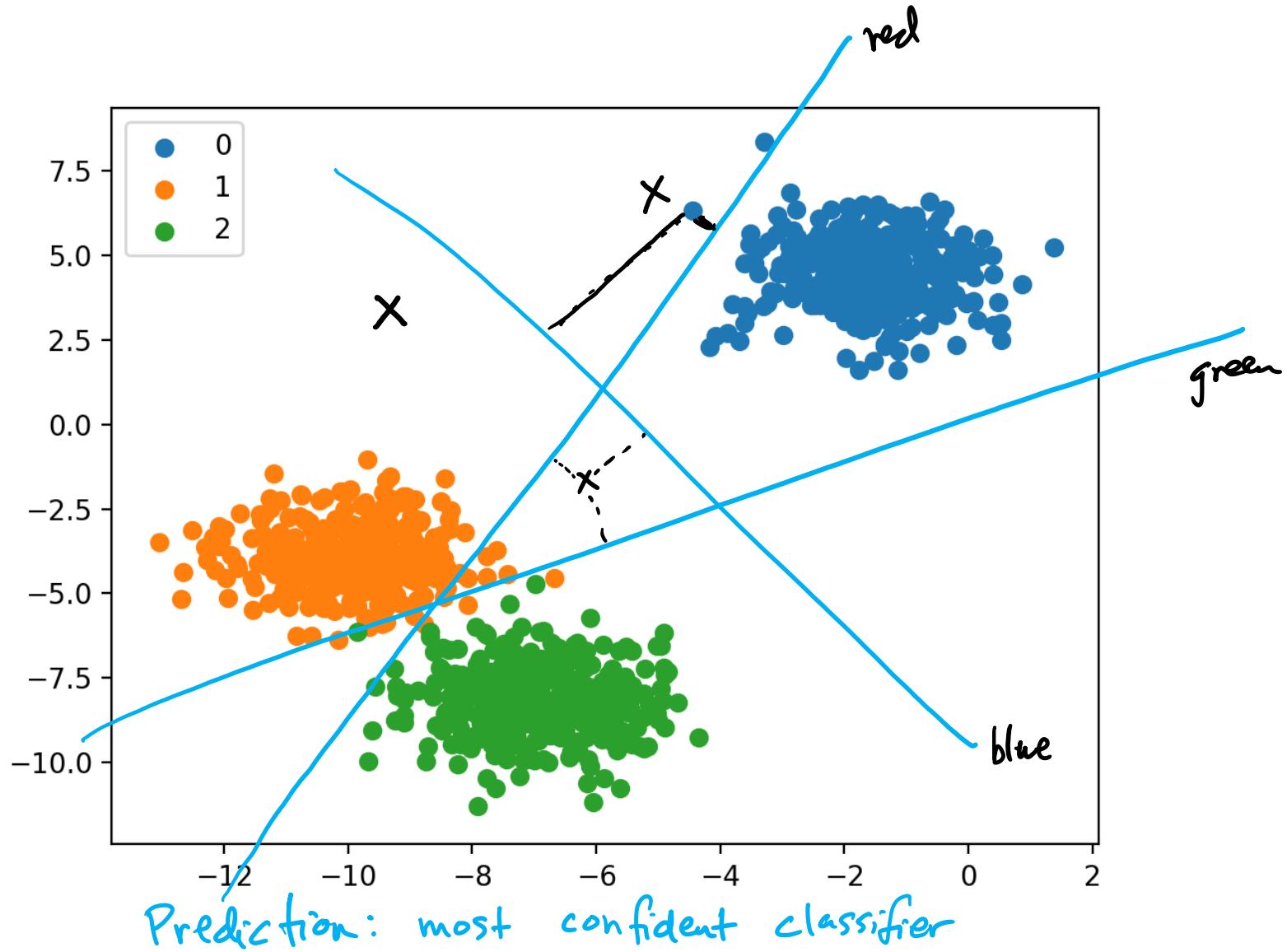
Multiclass Classification

- One vs. the rest
- All pairs

} reduce to a set of binary problems

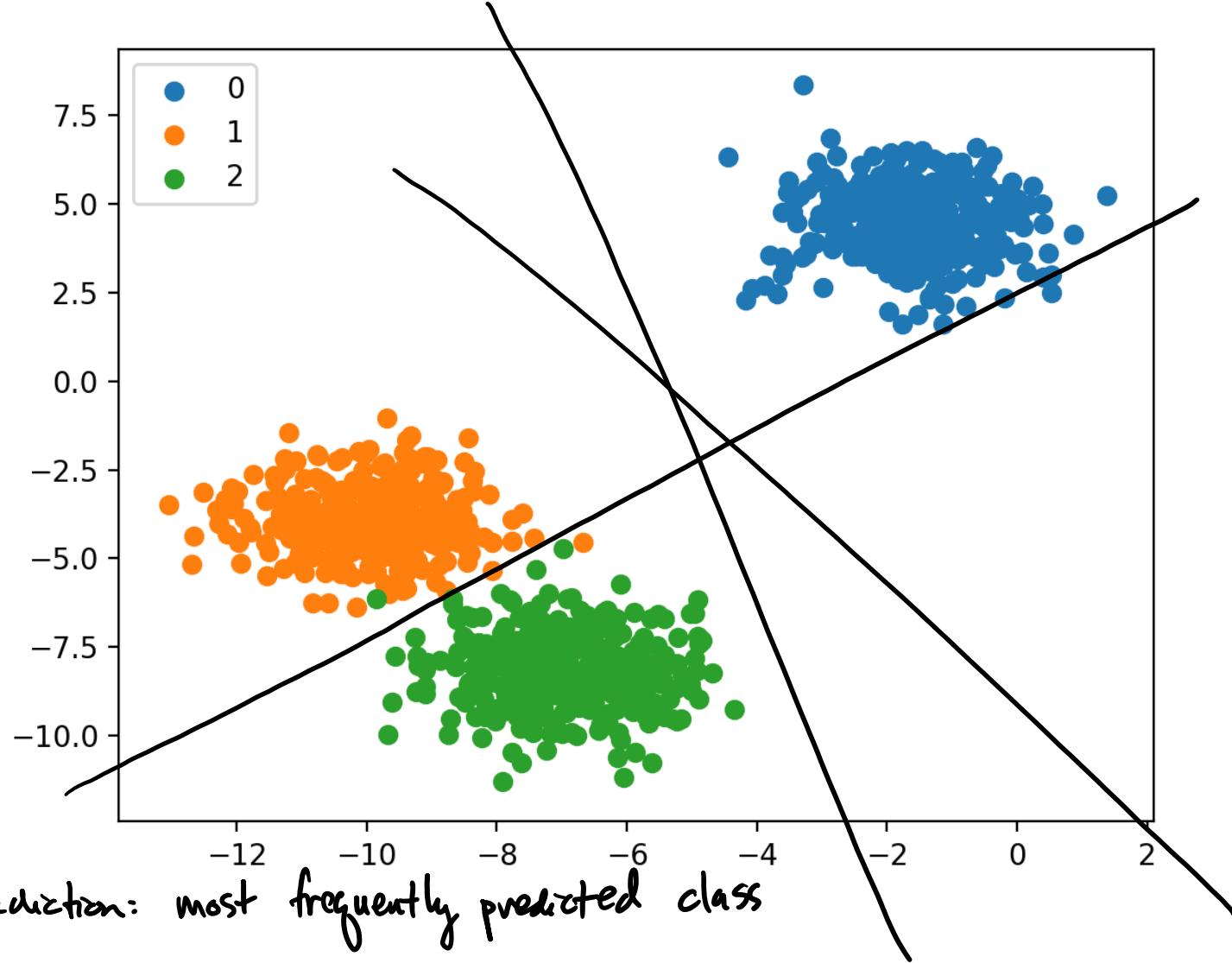
- Multiclass max margin and soft margin linear classifiers
- Multiclass loss functions
 - multiclass hinge losses
 - multiclass logistic loss

One vs. The Rest



$$\binom{K}{2}$$

All Pairs



Multiclass Linear Classification

- K classes
- For each $k \in \{1, \dots, K\}$, there is an affine function

$$f_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + b_k \quad \text{score function}$$

- The final classifier is

$$\mathbf{x} \mapsto \operatorname{argmax}_k f_k(\mathbf{x})$$

- We have already seen an example: **LDA**
- How can we extend other linear classification methods to learn these affine functions?

Max-Margin Formulation

- Training data $(x_1, y_1), \dots, (x_n, y_n)$ $y_i \in \{1, \dots, K\}$
- The training data is *linearly separable* iff there exist $w_1, b_1, \dots, w_K, b_K$ such that

$$\underbrace{w_{y_i}^T x_i + b_{y_i}}_{\text{score function for class } y_i} > \underbrace{w_k^T x_i + b_k}_{\text{score function for class } k} \quad \forall i \quad \forall k \neq y_i$$

- Max-margin separating hyperplane

$$\min_{\{w_k, b_k\}} \frac{1}{2} \sum_{k=1}^K \|w_k\|^2$$

$$\text{s.t. } (w_{y_i}^T x_i + b_{y_i}) - (w_k^T x_i + b_k) \geq 1 \quad \forall i \quad \forall k \neq y_i$$

- Reduces to what we have studied when $K = 2$

Soft-Margin Formulation

- Crammer-Singer version:

$$\begin{aligned} \min_{\{w_k, b_k\}, \{\xi_i\}} \quad & \frac{1}{2} \sum_k \|w_k\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & (w_{y_i}^T x_i + b_{y_i}) - (w_k^T x_i + b_k) \geq 1 - \xi_i \quad \forall i, k \neq y_i \\ & \xi_i \geq 0 \quad \forall i \end{aligned}$$

- Weston-Watkins version:

$$\begin{aligned} \min_{\{w_k, b_k\}, \{\xi_{ik}\}} \quad & \frac{1}{2} \sum \|w_k\|^2 + \frac{C}{n} \sum_{i=1}^n \sum_{k \neq y_i} \xi_{ik} \\ \text{s.t.} \quad & (w_{y_i}^T x_i + b_{y_i}) - (w_k^T x_i + b_k) \geq 1 - \xi_{ik} \quad \forall i, k \neq y_i \\ & \xi_{ik} \geq 0 \end{aligned}$$

Multiclass Hinge Losses

- ERM formulation: See my lecture notes for Crammer-Singer and Weston-Watkins
- These multiclass methods can be kernelized

$$\min_{w_1, b_1, \dots, w_k, b_k} \frac{1}{n} \sum L(y_i, f(x_i)) + \Omega(f)$$

$$f(x) = (w_1^T x + b_1, \dots, w_k^T x + b_k)$$

Multiclass Logistic Regression

- Recall

$$\eta_k(\mathbf{x}) = \Pr(Y = k \mid \mathbf{X} = \mathbf{x}) \quad \text{Hx} \quad \sum_{k=1}^K \eta_k(\mathbf{x}) = 1$$

- Define the *softmax* function, which maps

$$\mathbb{R}^K \rightarrow \{(p_1, \dots, p_K) \mid p_k \geq 0, \sum_k p_k = 1\}$$

probability simplex

by $\psi(\mathbf{v}) = (\psi_1(\mathbf{v}), \dots, \psi_K(\mathbf{v}))$ where

$$\psi_k(\mathbf{v}) := \frac{\exp(v_k)}{\sum_{l=1}^K \exp(v_l)}$$

- The *multinomial logistic regression* model assumes the existence of $\mathbf{w}_1, b_1, \dots, \mathbf{w}_K, b_K$ such that

$$\eta_k(\mathbf{x}) = \text{Hx} \quad \eta_k(w_1^\top \mathbf{x} + b_1, \dots, w_K^\top \mathbf{x} + b_K)$$

where $f_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + b_k$.

MLE/ERM

- Let θ collect $w_1, b_1, \dots, w_K, b_K$.
- Denote

$$\begin{aligned} f_{\theta}(x) &= (f_1(x), \dots, f_K(x)) \\ &= (w_1^T x + b_1, \dots, w_K^T x + b_K) \end{aligned}$$

- The maximum likelihood estimation problem can be written

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n L(y_i, f_{\theta}(x_i))$$

where

$$L(y, v) = -\ln \varphi_y(v)$$

is the *multinomial logistic loss*.

- This loss is convex as a function of v
- Agrees with binary version when $K = 2$

Connection to Cross-Entropy

- If $\mathbf{p} = (p_1, \dots, p_K)$ and $\mathbf{q} = (q_1, \dots, q_K)$ are two probability vectors, the *cross-entropy* between them is

$$CE(\mathbf{q}, \mathbf{p}) := - \sum_k q_k \ln(p_k)$$

$$\mathbf{e}_k = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ | \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{kth position}$$

- Let \mathbf{e}_k denote the k th standard basis vector in \mathbb{R}^K
- Then

$$L(y, v) = CE(e_y, \psi(v)) \quad \text{"cross entropy loss"}$$

Hence the cross-entropy loss and multinomial logistic regression loss are the same.

- Preview: If θ contains the weights of a neural network $f_\theta(x)$ with K nodes in the output layer, the network is trained by solving

$$\min_{\theta} \sum_{i=1}^n CE(e_{y_i}, \psi(f_\theta(x))) = \sum_{i=1}^n L(y_i, f_\theta(x_i))$$