# Multiclass linear prediction

DSE 220
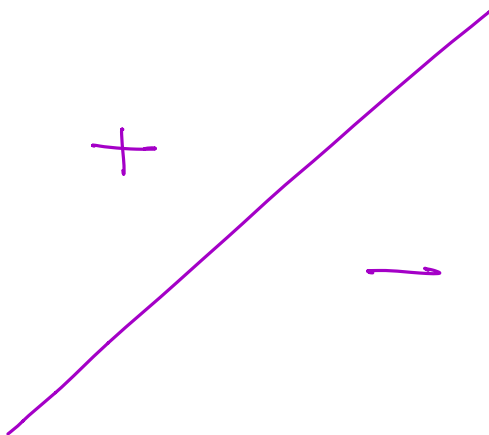
# Topics we'll cover

1. Multiclass logistic regression

2. Multiclass Perceptron

3. Multiclass support vector machines

# Multiclass classification

Of the classification methods we have studied so far, which seem inherently binary?

- Nearest neighbor?

- Generative models?

- Linear classifiers?

# The main idea

Remember Gaussian generative models...

Each class $j$ has an associated function

$$F_j(x) = \ln\left(\pi_j \; P_j(x)\right)$$

weight of class $j$ — Gaussian density for class $j$

To classify a new point $x$:

- Evaluate $F_1(x), F_2(x), .., F_k(x)$    [if $k$ classes]
- The biggest value wins $\rightarrow$ we predict that class

Gaussian case:
- $F_j(x)$ quadratic
- But linear if covariance matrices all equal

For linear classification, each class $(j = 1, 2, .., k)$ gets its own linear function.

Class 1: $\qquad w_1 \cdot x + b_1$

Class 2: $\qquad w_2 \cdot x + b_2$

$\qquad\qquad\qquad \vdots$

Class k: $\qquad w_k \cdot x + b_k$

if $x \in \mathbb{R}^d$ then

$w_1, .., w_k \in \mathbb{R}^d$

and

$b_1, .., b_k \in \mathbb{R}$

To predict the class of a new point $x$:

∘ Evaluate all $k$ of these functions

∘ Largest value wins

Linear classification in the multiclass setting

# From binary to multiclass logistic regression

**Binary** logistic regression: for $\mathcal{X} = \mathbb{R}^d$, classifier given by $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$:

$$\Pr(y = 1|x) = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}}$$

# From binary to multiclass logistic regression

**Binary** logistic regression: for $\mathcal{X} = \mathbb{R}^d$, classifier given by $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$:

$$\Pr(y = 1|x) = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}}$$

Labels $\mathcal{Y} = \{1, 2, \ldots, k\}$: specify a classifier by $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k \in \mathbb{R}$:

$$\Pr(y = j|x) \propto e^{w_j \cdot x + b_j}$$

# From binary to multiclass logistic regression

**Binary** logistic regression: for $\mathcal{X} = \mathbb{R}^d$, classifier given by $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$:

$$\Pr(y = 1|x) = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}}$$

Labels $\mathcal{Y} = \{1, 2, \ldots, k\}$: specify a classifier by $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k \in \mathbb{R}$:

$$\Pr(y = j|x) \;\propto\; e^{w_j \cdot x + b_j}$$

$\underset{\uparrow \text{ proportional}}{=}$

- What is the fully normalized form of the probability?

$$\Pr(y = j \,|\, x) \;=\; \frac{e^{w_j \cdot x + b_j}}{\sum\limits_{\ell=1}^{k} e^{w_\ell \cdot x + b_\ell}} \qquad \text{"softmax"}$$

- Given a point $x$, which label to predict?

$$\underset{j}{\arg\max} \quad w_j \cdot x + b_j$$

# Multiclass logistic regression

- **Label space**: $\mathcal{Y} = \{1, 2, \ldots, k\}$
- **Parametrized classifier**: $w_1, \ldots, w_k \in \mathbb{R}^d$, $b_1, \ldots, b_k \in \mathbb{R}$:

$$\Pr(y = j | x) = \frac{e^{w_j \cdot x + b_j}}{e^{w_1 \cdot x + b_1} + \cdots + e^{w_k \cdot x + b_k}}$$

- **Prediction**: given a point $x$, predict label $\arg\max_j (w_j \cdot x + b_j)$.
- **Learning**: Given: $(x^{(1)}, y^{(1)}), \ldots, (x^{(n)}, y^{(n)})$.    $x^{(i)} \in \mathbb{R}^d$,   $y^{(i)} \in \{1, 2, \ldots, k\}$
  Find: $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k$ that maximize the likelihood

$$\prod_{i=1}^{n} \Pr(y^{(i)} | x^{(i)})$$

Taking negative log gives a convex minimization problem.

# Multiclass Perceptron

Setting: $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{1, 2, \ldots, k\}$

**Model:** $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k \in \mathbb{R}$

**Prediction:** On instance $x$, predict label $\arg\max_j (w_j \cdot x + b_j)$

# Multiclass Perceptron

Setting: $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{1, 2, \ldots, k\}$

**Model:** $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k \in \mathbb{R}$

**Prediction:** On instance $x$, predict label $\arg\max_j(w_j \cdot x + b_j)$

**Learning.** Given training set $(x^{(1)}, y^{(1)}), \ldots, (x^{(n)}, y^{(n)})$:

- Initialize $w_1 = \cdots = w_k = 0$ and $b_1 = \cdots = b_k = 0$
- Repeat while some training point $(x, y)$ is misclassified:

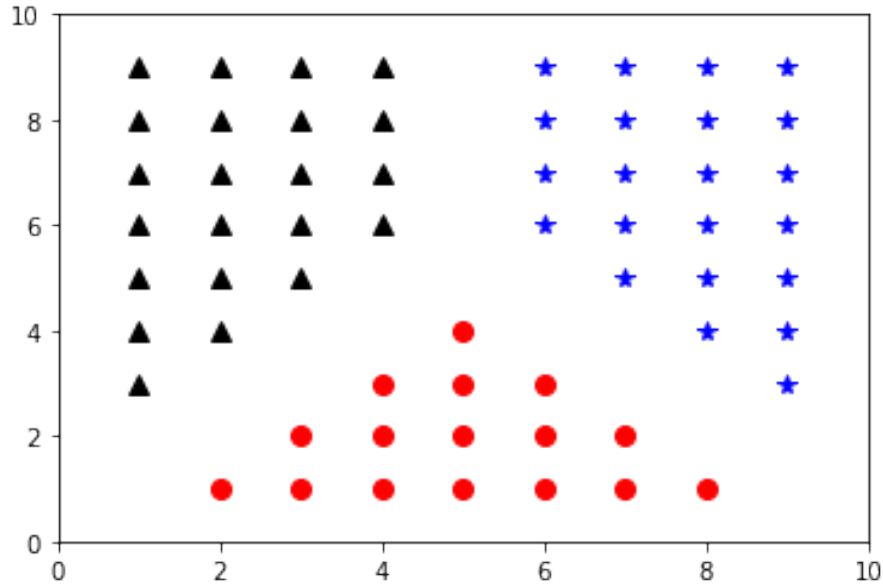  for correct label $y$:    $w_y = w_y + x$
  $$b_y = b_y + 1$$

  for predicted label $\widehat{y}$:    $w_{\widehat{y}} = w_{\widehat{y}} - x$
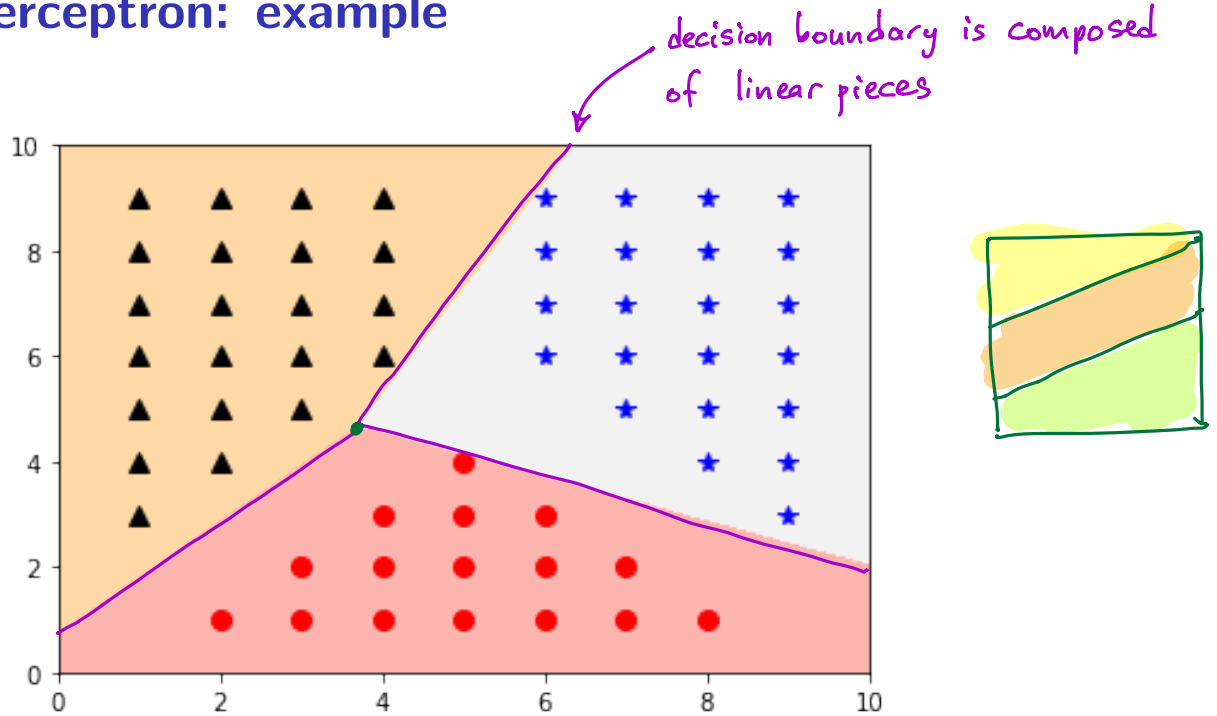  $$b_{\widehat{y}} = b_{\widehat{y}} - 1$$

only
update
two
classes

$\widehat{y}$ is
$\arg\max_j w_j \cdot x + b_j$

# Multiclass Perceptron: example

# Multiclass Perceptron: example

decision boundary is composed
of linear pieces

# Multiclass SVM

**Model:** $w_1, \ldots, w_k \in \mathbb{R}^d$ and $b_1, \ldots, b_k \in \mathbb{R}$

**Prediction:** On instance $x$, predict label $\arg\max_j (w_j \cdot x + b_j)$

**Learning.** Given training set $(x^{(1)}, y^{(1)}), \ldots, (x^{(n)}, y^{(n)})$:

*maximize margins*

*penalty for using slack*

*Convex program*

$$\min_{w_1, \ldots, w_k \in \mathbb{R}^d, b_1, \ldots, b_k \in \mathbb{R}, \xi \in \mathbb{R}^n} \sum_{j=1}^{k} \|w_j\|^2 + C \sum_{i=1}^{n} \xi_i$$

$$w_{y^{(i)}} \cdot x^{(i)} + b_{y^{(i)}} - w_y \cdot x^{(i)} - b_y \geq 1 - \xi_i \quad \text{for all } i, \text{ all } y \neq y^{(i)} \quad ?$$
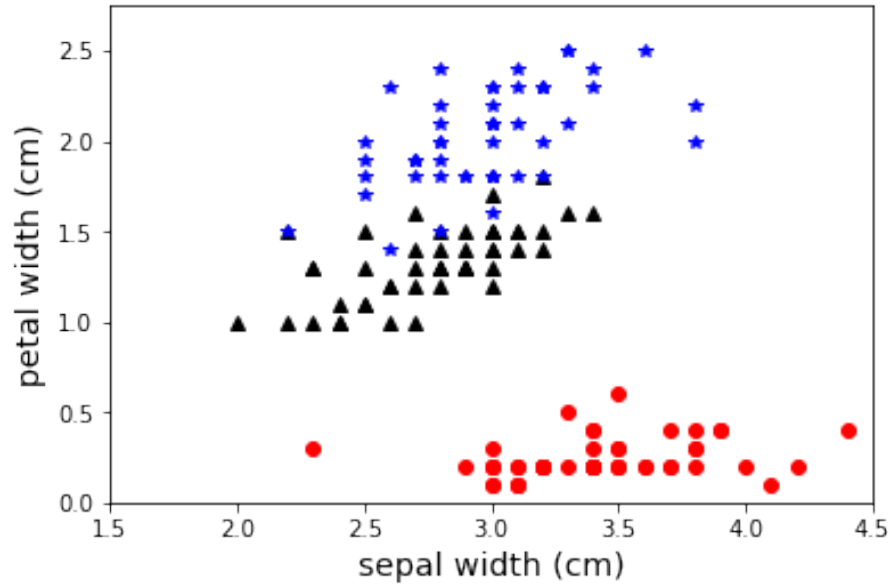
$$\xi \geq 0$$
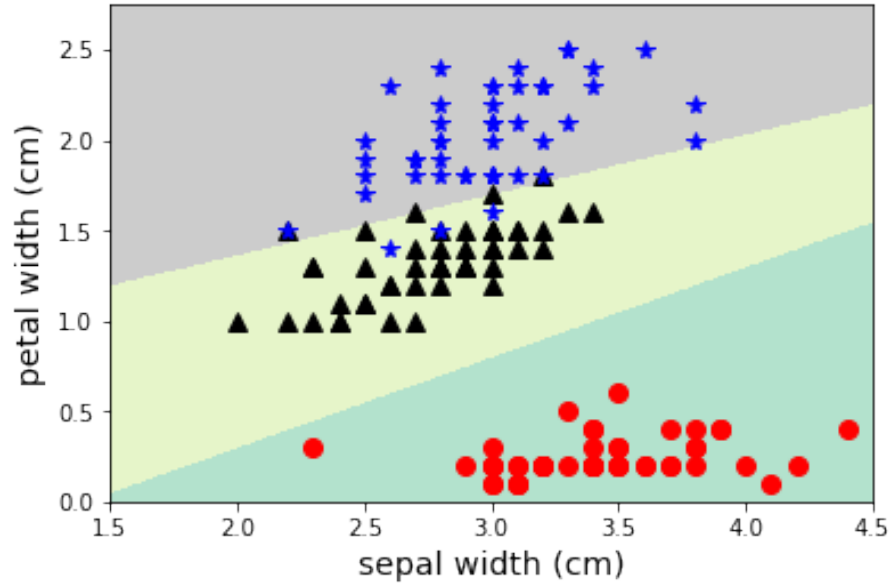
Point $x^{(i)}$ has true label $y^{(i)}$

$\therefore$ Want $\quad w_{y^{(i)}} \cdot x^{(i)} + b_{y^{(i)}} > w_y \cdot x^{(i)} + b_y \quad$ for all other labels $y \neq y^{(i)}$

As before, replace with $\quad w_{y^{(i)}} \cdot x^{(i)} + b_{y^{(i)}} \geq w_y \cdot x^{(i)} + b_y + 1 - \xi_i \quad$ *slack*

# Multiclass SVM example: iris

# Multiclass SVM example: iris

# Multiclass SVM

Given training set $(x^{(1)}, y^{(1)}), \ldots, (x^{(n)}, y^{(n)})$:

$$\min_{w_1,\ldots,w_k \in \mathbb{R}^d, b_1,\ldots,b_k \in \mathbb{R}, \xi \in \mathbb{R}^n} \sum_{j=1}^{k} \|w_j\|^2 + C \sum_{i=1}^{n} \xi_i$$

$$w_{y^{(i)}} \cdot x^{(i)} + b_{y^{(i)}} - w_y \cdot x^{(i)} - b_y \geq 1 - \xi_i \quad \text{for all } i, \text{ all } y \neq y^{(i)}$$

$$\xi \geq 0$$

Once again, a convex optimization problem.

Question: how many variables and constraints do we have?

**Variables**

① The classifiers: $k(d+1)$

② Slack variables: $n$

**Constraints**

For each data pt: want the correct label to beat the remaining $k-1$ labels.

∴ Total of $n(k-1)$ constraints.

Back to binary setting:

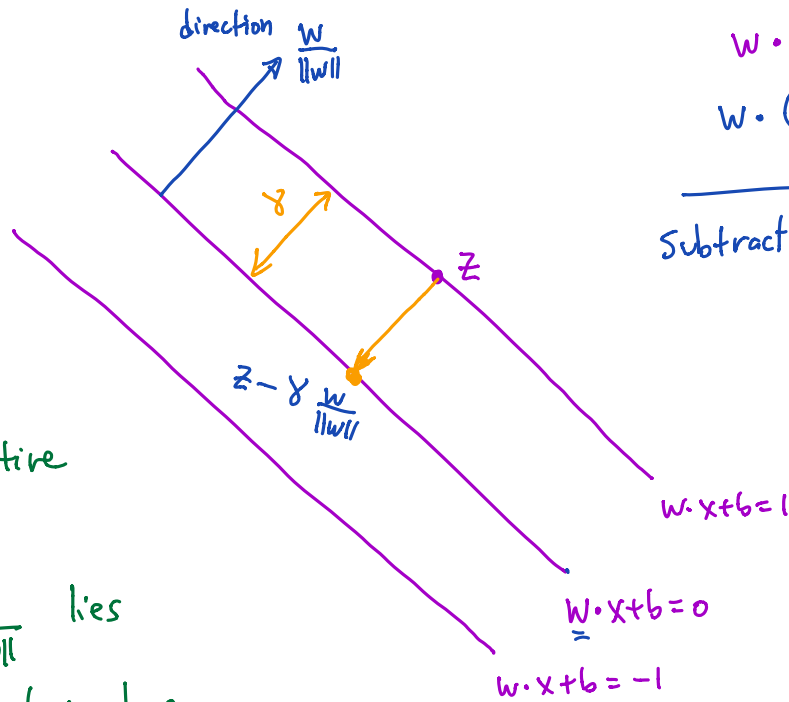$$y^{(i)} (w \cdot x^{(i)} + b) > 0 \qquad \text{for all} \quad i = 1..n$$

multiply $w, b$ by a large enough constant

$$\left\{ y^{(i)} (w \cdot x^{(i)} + b) \geq 1 \qquad \text{for all} \quad i = 1..n \right\}$$

if in this form, have a nice expression for the margin

$$\left( \text{margin} = \frac{1}{\|w\|} \right)$$

direction $\frac{w}{\|w\|}$

$$w \cdot z + b = 1$$
$$w \cdot \left( z - \gamma \frac{w}{\|w\|} \right) + b = 0$$

Subtract: $\gamma \|w\| = 1$
$$\Rightarrow \gamma = \frac{1}{\|w\|}$$

$\gamma$

$z$

$z - \gamma \frac{w}{\|w\|}$

$w \cdot x + b = 1$

$w \cdot x + b = 0$

$w \cdot x + b = -1$

① Let $z$ be <u>any</u> point on the positive boundary.

② Then $z - \gamma \frac{w}{\|w\|}$ lies on the decision boundary.

Worksheet 10
#1