

# Unit 5

## Linear Classification & Logistic Regression

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 5307: Introduction to Machine Learning, Sp23

# Learning objectives

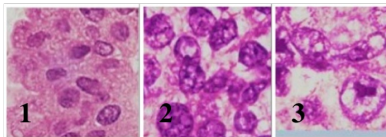
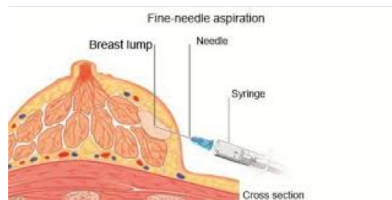
- Understand classification problems in machine learning:
  - Identify features, labels; binary vs multiclass; linear vs nonlinear
  - visualize scatterplots and decision regions
- For binary classification problems, understand ...
  - linear classifiers, separating hyperplanes, linearly separable data
  - effect of feature transformations
  - why LS linear regression doesn't work well
  - logistic regression: logistic function, cross-entropy loss, ML fitting, regularization
  - common error metrics: accuracy, precision, recall, F1
  - the effect of the decision threshold, ROC, AUC
- For multiclass classification problems, understand ...
  - solutions that use multiple binary classifiers
  - multinomial logistic regression: softmax function, cross-entropy loss, ML fitting
- How to implement and assess classification using **sklearn**

# Outline

- Motivating Example: Diagnosing Breast Cancer
- Binary Classification
- Binary Logistic Regression
- Multiclass Classification
- Multinomial Logistic Regression
- Measuring Accuracy in Classification

# Diagnosing Breast Cancer

- **Fine-needle aspiration** of suspicious breast lumps:
  - Tissue is stained & viewed under microscope
  - Cytopathologist visually inspects cells and takes several measurements
  - Also tries to diagnose as **benign** or **malignant**
- Would like to improve diagnosis
- Can machine-learning help?



Grades of carcinoma cells

<http://breast-cancer.ca/5a-types/>

# The Wisconsin Breast Cancer Data Set

Univ. of Wisconsin study:

- 683 samples
- 9 features (on right)
- target: malignant or benign
  - ground-truth was assessed using a biopsy

#	Attribute	Domain
-----		
1.	Sample code number	id number
2.	Clump Thickness	1 - 10
3.	Uniformity of Cell Size	1 - 10
4.	Uniformity of Cell Shape	1 - 10
5.	Marginal Adhesion	1 - 10
6.	Single Epithelial Cell Size	1 - 10
7.	Bare Nuclei	1 - 10
8.	Bland Chromatin	1 - 10
9.	Normal Nucleoli	1 - 10
10.	Mitoses	1 - 10
11.	Class:	(2 for benign, 4 for malignant)

Data:

<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data>

Explanation:

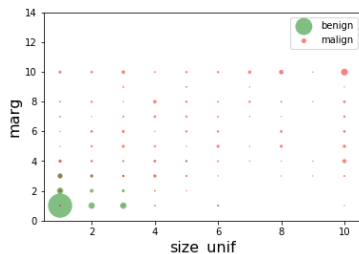
<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.names>

Original paper:

O.L. Mangasarian, W.N. Street, and W.H. Wolberg, "Breast cancer diagnosis and prognosis via linear programming," *Operations Research*, 1995.

# Visualizing the data

- Choose two features to visualize
- Plot target count vs. features using variable-radius dots (all are discrete)



- How would we classify a test feature  $x = [x_1, x_2]$ ?

```
# Compute the bin edges for the 2d histogram
x1val = np.array(list(set(X[:,0]))).astype(float)
x2val = np.array(list(set(X[:,1]))).astype(float)
x1, x2 = np.meshgrid(x1val,x2val)
x1e= np.hstack((x1val,np.max(x1val)+1))
x2e= np.hstack((x2val,np.max(x2val)+1))

# Make a plot for each class
yval = list(set(y))
color = ['g','r']
for i in range(len(yval)):
    I = np.where(y==yval[i])[0]
    cnt, x1e, x2e = np.histogram2d(X[I,0],X[I,1],[x1e,x2e])
    x1, x2 = np.meshgrid(x1val,x2val)
    plt.scatter(x1.ravel(), x2.ravel(), s=2*cnt.ravel(),alpha=0.5,
               c=color[i],edgecolors='none')

plt.ylim([0,14])
plt.legend(['benign','malign'], loc='upper right')
plt.xlabel(xnames[0], fontsize=16)
plt.ylabel(xnames[1], fontsize=16)
```

# Outline

- Motivating Example: Diagnosing Breast Cancer
- **Binary Classification**
- Binary Logistic Regression
- Multiclass Classification
- Multinomial Logistic Regression
- Measuring Accuracy in Classification

# Binary classification

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$ , our goal is to **classify** a test vector  $\mathbf{x}$  as  $y \in \{-1, 1\}$  (one of two “**classes**”)

- Unlike regression, the target  $y$  is now **binary**
- Could also use  $\{0, 1\}$  or other target labels

- Many applications:

- Are these cells cancerous or not?
- Is this a weed or a crop?

- Mathematically, want to design a **classifier**

$$f(\mathbf{x}) = \hat{y} \in \{-1, 1\}$$

such that  $\hat{y} = y$  with high probability

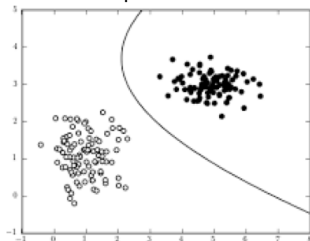
- Note:  $f(\cdot)$  is defined by its **decision regions**:

$$\mathcal{R}_{-1} \triangleq \{\mathbf{x} : f(\mathbf{x}) = -1\} \text{ and } \mathcal{R}_1 \triangleq \{\mathbf{x} : f(\mathbf{x}) = 1\}$$

Is this a weed or a crop?



Example with  $d = 2$ :





# Binary linear classification

- One option is **binary linear classification**:

1) compute the “**score**” or “**discriminant**”,  $z = b + \sum_{j=1}^d x_j w_j$

- $z$  is linear in the parameters  $b$  and  $w_j$
- $z$  is linear in the features  $x_j$

2) threshold the score to obtain  $\hat{y} = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$

The weights  $\mathbf{w} \triangleq [w_1, \dots, w_d]^T$  and intercept  $b$  are learned from the training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

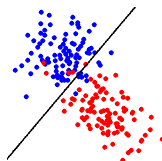
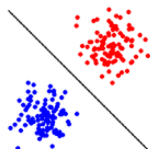
- Note that, at the **decision boundary**, we have

$$0 = z = b + \mathbf{x}^T \mathbf{w}.$$

Thus, the **hyperplane**

$$\{\mathbf{x} : b + \mathbf{x}^T \mathbf{w} = 0\}$$

separates the decision regions



- This **decision boundary is linear** in  $\mathbf{x}$  (see figures). When does this perform well?

# Linear separability

- Linear classification tends to perform well when the training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  is “linearly separable”

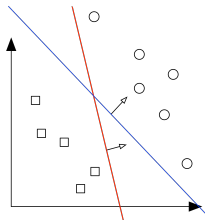
- Linearly separable means that there exists a hyperplane

$$\{\mathbf{x} : b + \mathbf{x}^T \mathbf{w} = 0\}$$

(for some  $b$  and  $\mathbf{w}$ ) that separates the samples  $\mathbf{x}_i$  according to their class  $y_i \in \{-1, 1\}$

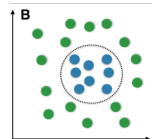
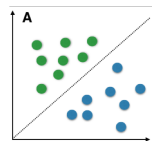
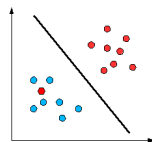
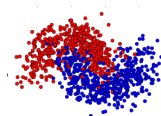
- On one side of hyperplane lies all  $\mathbf{x}_i$  for which  $y_i = 1$ , while on other other side lies all  $\mathbf{x}_i$  for which  $y_i = -1$
- Note: When such a separating hyperplane exists, it is usually non-unique

linear separability:



# Linear versus nonlinear classification

- Most datasets are *not* linearly separable!
  - There are many possible reasons why
  - See examples on right (except Fig. A)
- Still, linear classification is worth considering
  - It is relatively easy to understand
  - It facilitates feature selection (i.e., provides intuition about which features matter most)
  - It can incorporate nonlinear feature transformations
    - Fig. A: boundary is linear in  $(x_1, x_2)$
    - Fig. B: boundary is nonlinear in  $(x_1, x_2)$  but linear in a new feature  $x' \triangleq \sqrt{x_1^2 + x_2^2}$
  - It is often used as a building block (e.g., for neural networks, decision trees, etc.)



# Linear classification vs. LS linear regression

- Suppose we want to do *linear* classification. How exactly do we fit  $(b, w)$ ?
- Can we just **use LS**, as we have been doing?
  - In other words, choose  $(b, w)$  to minimize  $\text{RSS} = \left\| \mathbf{y} - \mathbf{A} \begin{bmatrix} b \\ w \end{bmatrix} \right\|^2$  and then output

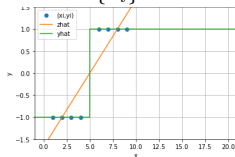
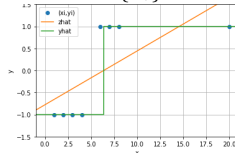
$$\hat{y} = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}, \quad \text{where } z = [1 \ x^T] \begin{bmatrix} b_{\text{ls}} \\ w_{\text{ls}} \end{bmatrix}$$

- **No, this can fail terribly!**

- Consider simple case of  $d = 1$  feature
- When  $\{x_i\}$  is “balanced” (top figure), works okay
- But when  $\{x_i\}$  is “imbalanced” (bottom), the least-squares regression line gets pulled to one side, and the threshold at  $z = 0$  makes errors

- What's the problem?

- **RSS is not the right loss function** for classification!  
(More details soon)

balanced  $\{x_i\}$ :imbalanced  $\{x_i\}$ :

# LS linear regression fails on breast-cancer classification!

- Let's try the same **LS approach** on the breast-cancer data (after converting targets to  $y \in \pm 1$ )

- Again, visualize decision boundary

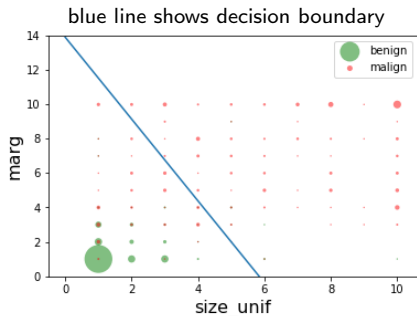
$$\{x : b + x^T w = 0\}$$

which, in this  $d = 2$  case, becomes

$$\{(x_1, x_2) : x_2 = -\frac{b}{w_2} - \frac{w_1}{w_2}x_1\}$$

- Because the  $\{x_i\}$  are **imbalanced**, the decision boundary is pulled north-east, and many red  $x_i$  are misclassified!

- Again we see that designing  $(b, w)$  to minimize RSS does not work well for linear classification!



```
yhat = regr.predict(X)
yhati = (yhat >= 0.5).astype(int)
acc = np.mean(yhati == y)
print("Accuracy on training data using two features = %f" % acc)
```

Accuracy on training data using two features = 0.922401

# Outline

- Motivating Example: Diagnosing Breast Cancer
- Binary Classification
- **Binary Logistic Regression**
- Multiclass Classification
- Multinomial Logistic Regression
- Measuring Accuracy in Classification

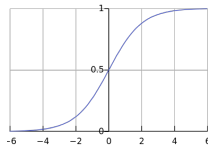
# Binary logistic regression

- **Linear classification** computes  $z = b + \mathbf{x}^\top \mathbf{w}$  and sets  $\hat{y} = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases}$
- How do we design the parameters  $(b, \mathbf{w})$ ?
  - We saw that minimizing RSS does not work well
- Idea: Given the score  $z$ , model the true label  $y \in \pm 1$  as a **random variable**
- The most popular version of this uses the probabilistic model

$$\Pr\{y=1 \mid z\} = \frac{e^z}{1 + e^z}, \quad \Pr\{y=-1 \mid z\} = \frac{1}{1 + e^z}$$

Note  $\Pr\{y=-1 \mid z\} + \Pr\{y=1 \mid z\} = 1 \quad \forall z$ , as required for a valid pmf.

- The larger that  $z$  is, the more likely that  $y = 1$
- When  $z = 0$ , it's equally likely that  $y = 1$  or  $y = -1$
- $\frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$  is known as the **“logistic function”** or the **“sigmoid”**



# Understanding the logistic model

Previously we considered  $\Pr\{y=1 \mid z\}$ . What about  $\Pr\{y=1 \mid x\}$ ?

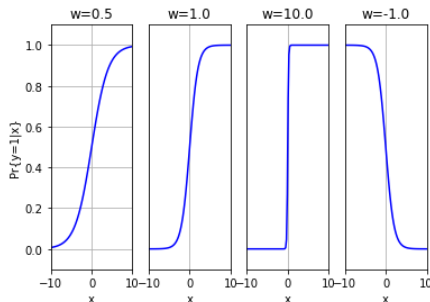
- Consider the simple case of a single scalar feature  $x$  and intercept  $b=0$ :

$$\Pr\{y=1 \mid x\} = \frac{1}{1 + e^{-z}} \text{ for } z = wx$$

- Sign of  $w$  determines increasing vs. decreasing in  $x$
- Magnitude of  $w$  determines sharpness of  $x$ -domain transition
- Transition centered at  $x = 0$
- For a more general  $b$ , we have

$$z = b + wx = w\left(\frac{b}{w} + x\right).$$

Thus transition is centered at  $x = -\frac{b}{w}$



$\Pr\{y=1 \mid x\}$  versus  $x$  for various  $w$



# Maximum likelihood estimation

- Given training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , we can fit the model parameters  $(b, \mathbf{w})$  using **maximum likelihood (ML) estimation**:

## ML Estimation

- 1) Define a **likelihood function**  $p(\mathbf{y}|\mathbf{X}, [\begin{smallmatrix} b \\ \mathbf{w} \end{smallmatrix}])$  with model parameters  $[\begin{smallmatrix} b \\ \mathbf{w} \end{smallmatrix}]$ 
    - As usual,  $\mathbf{y} \triangleq [y_1, \dots, y_n]^\top$  and  $\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$
  - 2) The ML model parameters are
 
$$\begin{aligned} (b_{\text{ml}}, \mathbf{w}_{\text{ml}}) &\triangleq \arg \max_{b, \mathbf{w}} p(\mathbf{y}|\mathbf{X}, [\begin{smallmatrix} b \\ \mathbf{w} \end{smallmatrix}]) \\ &= \arg \max_{b, \mathbf{w}} \ln p(\mathbf{y}|\mathbf{X}, [\begin{smallmatrix} b \\ \mathbf{w} \end{smallmatrix}]) \\ &= \arg \min_{b, \mathbf{w}} \{ -\ln p(\mathbf{y}|\mathbf{X}, [\begin{smallmatrix} b \\ \mathbf{w} \end{smallmatrix}]) \} \end{aligned}$$
- Recall that we previously applied ML estimation to linear regression:
    - There, the likelihood was  $p(\mathbf{y}|\mathbf{X}, \beta) = \mathcal{N}(\mathbf{y}; \mathbf{A}\beta, \sigma^2 \mathbf{I})$  with  $\mathbf{A} = [\mathbf{1} \ \mathbf{X}]$
    - $\Rightarrow -\ln p(\mathbf{y}|\mathbf{X}, \beta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\beta\|^2 + \text{constant}$
    - $\Rightarrow \beta_{\text{ml}} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{A}\beta\|^2 = \arg \min_{\beta} \text{RSS}(\beta) = \beta_{\text{ls}} \dots \text{the LS fit!}$
  - For logistic regression, we use ML estimation with a different likelihood

# ML estimation for logistic regression (LR)

- Logistic regression assumes that  $y_i$  depends on  $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$  as follows:

$$\Pr\{y_i=1 \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\} = \frac{e^{z_i}}{1+e^{z_i}}, \quad \Pr\{y_i=-1 \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\} = \frac{1}{1+e^{z_i}}, \quad z_i = b + \mathbf{x}_i^T \mathbf{w}$$

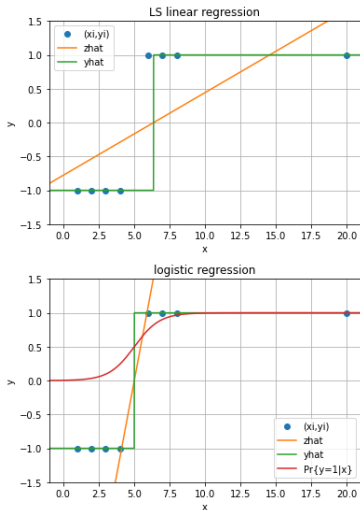
- Thus we have

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}) &= \prod_{i=1}^n p(y_i \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}) && \text{via independence} \\ -\ln p(\mathbf{y} \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}) &= -\sum_{i=1}^n \ln p(y_i \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}) && \text{since } \ln(cd) = \ln c + \ln d \\ &= -\sum_{i=1}^n \frac{y_i+1}{2} \ln \Pr\{y_i=1 \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\} + (1 - \frac{y_i+1}{2}) \ln \Pr\{y_i=-1 \mid \mathbf{X}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\} \\ &\quad \text{since } \frac{y_i+1}{2} = \begin{cases} 1 & y_i = 1 \\ 0 & y_i = -1 \end{cases} \dots \text{switches "on" or "off"} \\ &= -\sum_{i=1}^n \frac{y_i+1}{2} (z_i - \ln[1 + e^{z_i}]) + (1 - \frac{y_i+1}{2}) (0 - \ln[1 + e^{z_i}]) \\ &= -\sum_{i=1}^n \left( \frac{y_i+1}{2} z_i - \frac{y_i+1}{2} \ln[1 + e^{z_i}] - \ln[1 + e^{z_i}] + \frac{y_i+1}{2} \ln[1 + e^{z_i}] \right) \\ &= -\sum_{i=1}^n \left( \frac{y_i+1}{2} z_i - \ln[1 + e^{z_i}] \right) && \text{where } z_i = b + \mathbf{x}_i^T \mathbf{w} \end{aligned}$$

- The ML estimate of  $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$  is the one minimizing this negative log likelihood

# Comparison of LS vs. LR fits on a simple dataset

- Consider ML-fitting of the imbalanced 1D dataset that we saw earlier:
  - The top plot shows the LS regression line  $z_{ls}(x) \triangleq b_{ls} + w_{ls}x$ , as well as the binary prediction  $\hat{y}_{ls}(x) = \text{sgn}(z_{ls}(x))$
  - The bottom plot shows the LR score  $z_{lr}(x) \triangleq b_{lr} + w_{lr}x$ , as well as the binary prediction  $\hat{y}_{lr}(x) = \text{sgn}(z_{lr}(x))$
  - The plots show that LS is distracted by the outlier sample, while LR is not!
  - The bottom plot also shows  $\Pr\{y = 1 | x\}$  for  $b_{lr}$  and  $w_{lr}$ . As  $x$  increases, LR becomes more confident that  $y = 1$



# Binary cross-entropy loss

In summary...

- When  $y_i \in \{-1, 1\}$ , the ML weights for binary logistic regression are

$$(b_{\text{ml}}, \mathbf{w}_{\text{ml}}) \triangleq \arg \min_{b, \mathbf{w}} \sum_{i=1}^n \left( \ln[1 + e^{z_i}] - \frac{y_i + 1}{2} z_i \right) \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$

which must be solved numerically

- When  $y_i \in \{0, 1\}$ , the ML weights for binary logistic regression are

$$(b_{\text{ml}}, \mathbf{w}_{\text{ml}}) \triangleq \arg \min_{b, \mathbf{w}} \sum_{i=1}^n \left( \ln[1 + e^{z_i}] - y_i z_i \right) \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$

which is a popular alternative expression

- The summation is known as the “**logistic loss**” or “**binary cross-entropy loss**”

# Adding regularization

Assuming  $y_i \in \{0, 1\}$  for the following expressions ...

- Usually, **L2 regularization** is used with the binary cross-entropy loss:

$$(b_{lr}, \mathbf{w}_{lr}) = \arg \min_{b, \mathbf{w}} \left\{ \sum_{i=1}^n (\ln[1 + e^{z_i}] - y_i z_i) + \alpha \|\mathbf{w}\|^2 \right\} \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$

- When the training data is linearly separable,  $\|\mathbf{w}_{ml}\| = \infty$  without regularization!  
Using L2 regularization with small positive  $\alpha$  will keep  $\mathbf{w}_{lr}$  finite
  - With correlated features, larger  $\alpha$  may be helpful. Tune via cross-val
- Could instead use **L1 regularization**, and thus perform feature selection:
 
$$(b_{lr}, \mathbf{w}_{lr}) = \arg \min_{b, \mathbf{w}} \left\{ \sum_{i=1}^n (\ln[1 + e^{z_i}] - y_i z_i) + \alpha \|\mathbf{w}\|_1 \right\} \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$
  - As before, use L1 for feature-selection only. Once features are selected, use LR (with L2 & small  $\alpha$ ). Tune  $\alpha$  via cross-val performance of feature-selected LR
  - Computationally, L1 is much more expensive than L2

# Logistic regression in sklearn

In `sklearn`, there is a convenient `LogisticRegression` method:

- Note: L2 regularization is used by default
  - Don't forget to standardize  $X$ !
- The *inverse* regularization strength is controlled by the parameter  $C > 0$ 
  - So small regularization weight  $\alpha$  corresponds to *large*  $C$

```
logreg = linear_model.LogisticRegression(C=1e5)
```

```
Xs = preprocessing.scale(X)
logreg.fit(Xs, y)
```

```
LogisticRegression(C=100000.0, class_weight=None, dual=False,
                    fit_intercept=True, intercept_scaling=1, max_iter=100,
                    multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
                    solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
```

```
yhat = logreg.predict(Xs)
acc = np.mean(yhat == y)
print("Accuracy on training data = %f" % acc)
```

```
Accuracy on training data = 0.969253
```

# Outline

- Motivating Example: Diagnosing Breast Cancer
- Binary Classification
- Binary Logistic Regression
- **Multiclass Classification**
- Multinomial Logistic Regression
- Measuring Accuracy in Classification

# Multiclass classification

- What if there are  $K > 2$  classes?
  - Binary classification corresponds to  $K = 2$
- Goal: Given training data  $\{(x_i, y_i)\}_{i=1}^n$  with  $x_i \in \mathbb{R}^d$  and  $y_i \in \{1, \dots, K\}$ , we want to **classify** a test vector  $x$  as  $y \in \{1, \dots, K\}$ 
  - Unlike regression, the target  $y$  is **categorical**
- Mathematically, we want to design a **classifier**

$$f(x) = \hat{y} \in \{1, \dots, K\}$$

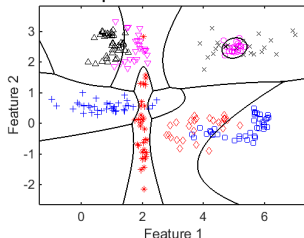
such that  $\hat{y} = y$  with high probability

- Again,  $f(\cdot)$  is defined by its **decision regions**:

$$\mathcal{R}_k \triangleq \{x : f(x) = k\} \quad \text{for } k = 1, \dots, K$$

- Important: Classification problems have **categorical** targets, while regression problems have **ordinal** targets
  - If the target is **discrete**, but ordinal, it's probably a regression problem

example where  $K = 10$ :





# Multiclass classification using binary classifiers

Multiclass classification is difficult. Can we tackle it using *binary* methods? Yes!

Two main approaches:

- **1-vs-rest** (also called **1-vs-all**)

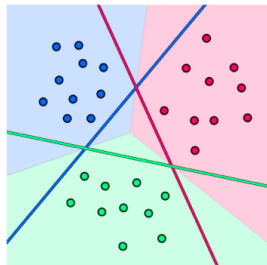
- for each class  $k = 1 \dots K$ , compute the score  $z_k$  that  $\mathbf{x}$  is **in- $k$  versus not-in- $k$**
- choose the class with highest score:  

$$\hat{k} = \arg \max_k z_k$$
- most common approach

- **1-vs-1:**

- for each class pair  $(k, l) | l \neq k$ , compute the score  $z_{kl}$  that  $\mathbf{x}$  is **in- $k$  versus in- $l$**
- choose the class with highest total score:  

$$\hat{k} = \arg \max_k \sum_{l \neq k} z_{kl}$$
- less common



# Outline

- Motivating Example: Diagnosing Breast Cancer
- Binary Classification
- Binary Logistic Regression
- Multiclass Classification
- **Multinomial Logistic Regression**
- Measuring Accuracy in Classification

# Multinomial logistic regression (MLR)

Multi-class extensions of binary logistic regression based on “1-vs-rest”:

- Test features  $\mathbf{x}$  are classified using a two-step approach:
  - For each class  $k = 1, \dots, K$ , compute a linear **score**  $z_k = b_k + \sum_{j=1}^d x_j w_{kj}$  that indicates whether  $\mathbf{x}$  is **in- $k$**  or **not-in- $k$** ,
  - Choose the winning class hypothesis as:  $\hat{y} = \arg \max_k z_k$
- Two options to design the coeffs  $\{(b_k, \mathbf{w}_k)\}_{k=1}^K$ , where  $\mathbf{w}_k \triangleq [w_{k1}, \dots, w_{kd}]^T$ :
  - 1 train  $K$  separate binary logistic regression models, or
  - 2 maximize a likelihood function constructed from the “**softmax**” model:

$$\Pr\{y=k \mid \mathbf{z}\} = \frac{e^{z_k}}{\sum_{l=1}^K e^{z_l}} \quad (\text{denominator ensures } \sum_{k=1}^K \Pr\{y=k \mid \mathbf{z}\} = 1 \ \forall \mathbf{z})$$

It's called “softmax” because, if  $z_{k_{\max}} \gg z_k$  for all  $k \neq k_{\max}$ , then

$$\frac{e^{z_k}}{\sum_{l=1}^K e^{z_l}} \approx \begin{cases} 1 & \text{if } k = k_{\max} \\ 0 & \text{if } k \neq k_{\max} \end{cases} = \text{one-hot coding of } k_{\max}$$

# One-hot label-coding for MLR

- Goal: Design  $\mathbf{b} = [b_1, \dots, b_K]^T$  &  $\mathbf{W} \triangleq [\mathbf{w}_1, \dots, \mathbf{w}_K]$  using **ML estimation**:

$$(\mathbf{b}_{\text{ml}}, \mathbf{W}_{\text{ml}}) \triangleq \arg \max_{\mathbf{b}, \mathbf{W}} p(\mathbf{y} | \mathbf{X}, \mathbf{b}, \mathbf{W}) = \arg \min_{\mathbf{b}, \mathbf{W}} \left\{ - \sum_{i=1}^n \ln p(y_i | \mathbf{X}, \mathbf{b}, \mathbf{W}) \right\}$$

- To help formulate this, we turn the categorical label  $y_i \in \{1, \dots, K\}$  into a binary vector  $\mathbf{y}_i \triangleq [y_{i1}, \dots, y_{iK}]^T$  using **one-hot-coding**, i.e.,

$$\mathbf{y}_i = [0 \dots 0, 1, 0 \dots 0]^T \quad \begin{matrix} \uparrow \\ y_i \text{th entry} \end{matrix} \Leftrightarrow y_{ik} = \begin{cases} 1 & \text{if } k = y_i \\ 0 & \text{if } k \neq y_i \end{cases} \quad \text{for all } i = 1, \dots, n$$

- Then, similar to the binary  $y_i \in \{0, 1\}$  case, we can write

$$\begin{aligned} -\ln p(y_i | \mathbf{X}, \mathbf{b}, \mathbf{W}) &= -\sum_{k=1}^K y_{ik} \ln \Pr\{y_i = k | \mathbf{X}, \mathbf{b}, \mathbf{W}\} \\ &= -\sum_{k=1}^K y_{ik} \ln \frac{e^{z_{ik}}}{\sum_l e^{z_{il}}} \quad \text{with } z_{ik} = b_k + \mathbf{x}_i^T \mathbf{w}_k \\ &= -\sum_{k=1}^K y_{ik} (z_{ik} - \ln [\sum_{l=1}^K e^{z_{il}}]) \\ &= \ln [\sum_{l=1}^K e^{z_{il}}] - \sum_{k=1}^K y_{ik} z_{ik}, \quad \text{since } \sum_{k=1}^K y_{ik} = 1 \quad \forall i \end{aligned}$$

# Multinomial cross-entropy loss

- Combining the results from the previous page, we get

$$(\mathbf{b}_{\text{ml}}, \mathbf{W}_{\text{ml}}) = \arg \min_{\mathbf{b}, \mathbf{W}} \underbrace{\left\{ \sum_{i=1}^n \left( \ln \left[ \sum_{k=1}^K e^{z_{ik}} \right] - \sum_{k=1}^K y_{ik} z_{ik} \right) \right\}}_{\triangleq J_{\text{mlr}}(\mathbf{b}, \mathbf{W})}, \quad z_{ik} = b_k + \mathbf{x}_i^T \mathbf{w}_k$$

where  $J_{\text{mlr}}(\mathbf{b}, \mathbf{W})$  is known as the **cross-entropy loss**

- In practice, we add  $\alpha \|\mathbf{W}\|_F^2$  (L2 regularization) or  $\alpha \|\mathbf{W}\|_1$  (L1 regularization), and tune  $\alpha$  via cross-validation
- With or without regularization, there is no closed-form expression for the optimal  $(\mathbf{b}, \mathbf{W})$ , but the optimization problems are convex and can be solved numerically. (More details given in Unit 6)

# Logistic Regression: Multinomial vs. Binary-OVR

- The `LogisticRegression` method in `sklearn` handles both binary and multiclass classification. It has many options:

```
LogisticRegression(C=100000.0, class_weight=None, dual=False,
                   fit_intercept=True, intercept_scaling=1, max_iter=100,
                   multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
                   solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
```

- With multiple classes, it is important to understand the `multi_class` option:
  - `multi_class = 'multinomial'`: This jointly trains  $(\mathbf{b}, \mathbf{W})$  according to the MLR approach described in the last few pages
  - `multi_class = 'ovr'`: For each  $k = 1 \dots K$ , this separately trains a one-vs-rest binary classifier  $(b_k, \mathbf{w}_k)$  using binary LR (i.e., option #1 on page 27)

The MLR version usually has better performance, but takes longer to train

# Outline

- Motivating Example: Diagnosing Breast Cancer
- Binary Classification
- Binary Logistic Regression
- Multiclass Classification
- Multinomial Logistic Regression
- Measuring Accuracy in Classification

# Performance metrics for binary classification

- In binary classification, there are 2 types of error:

- False Positive (or false alarm)
- False Negative (or missed detection)

contingency table or confusion matrix

	$\hat{y}=1$	$\hat{y}=0$
$y=1$	TP	FN
$y=0$	FP	TN

- The implications of these errors can be very different!

- e.g., consider breast cancer diagnosis

- Common machine-learning performance metrics:

- accuracy:**  $\Pr\{\hat{y} = y\} = \frac{TP+TN}{TP+FN+TN+FP}$

how often is the test correct?

- precision:**  $\Pr\{y=1 \mid \hat{y}=1\} = \frac{TP}{TP+FP}$

given a positive test, how often is the patient actually sick?

- recall:**  $\Pr\{\hat{y}=1 \mid y=1\} = \frac{TP}{TP+FN}$

given a sick patient, how often is the test correct?

- F1-score:**  $\left[ \frac{1}{2} \left( \frac{1}{\text{precision}} + \frac{1}{\text{recall}} \right) \right]^{-1}$

harmonic mean of precision &amp; recall

- Common metrics in medicine:

- sensitivity:**  $\Pr\{\hat{y}=1 \mid y=1\} = \frac{TP}{TP+FN}$

given a sick patient, how often is the test correct?

- specificity:**  $\Pr\{\hat{y}=0 \mid y=0\} = \frac{TN}{TN+FP}$

given a healthy patient, how often is the test correct?

[https://en.wikipedia.org/wiki/Evaluation\\_of\\_binary\\_classifiers](https://en.wikipedia.org/wiki/Evaluation_of_binary_classifiers)



# Breast cancer demo

- We now assess classification performance on the breast cancer demo
- Use 10-fold cross-validation
- `sklearn` includes support for computing **precision**, **recall**, **F1 score**, and **accuracy** (as defined on previous page)

```
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_validate

# Instantiate KFold object
nfold = 10
kf = KFold(n_splits=nfold, shuffle=True, random_state=0)

# Do cross-validation
scores = ('precision', 'recall', 'f1', 'accuracy')
cv_results = cross_validate(logreg, Xs, y, cv=kf, scoring=scores)

# Extract test metrics
prec = cv_results['test_precision']
rec = cv_results['test_recall']
f1 = cv_results['test_f1']
acc = cv_results['test_accuracy']

# Take average values of the metrics
precm = np.mean(prec)
recm = np.mean(rec)
f1m = np.mean(f1)
accm = np.mean(acc)

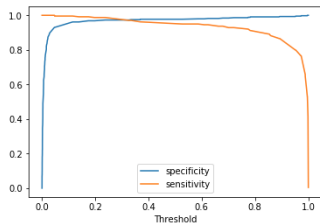
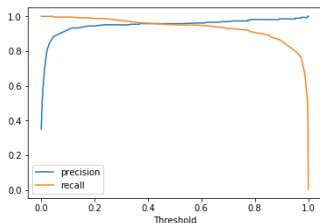
# Compute the standard errors
prec_se = np.std(prec, ddof=1)/np.sqrt(nfold)
rec_se = np.std(rec, ddof=1)/np.sqrt(nfold)
f1_se = np.std(f1, ddof=1)/np.sqrt(nfold)
acc_se = np.std(acc, ddof=1)/np.sqrt(nfold)

print('Precision = {0:.4f}, SE={1:.4f}'.format(precm, prec_se))
print('Recall = {0:.4f}, SE={1:.4f}'.format(recm, rec_se))
print('f1 = {0:.4f}, SE={1:.4f}'.format(f1m, f1_se))
print('Accuracy = {0:.4f}, SE={1:.4f}'.format(accm, acc_se))
```

```
Precision = 0.9432, SE=0.0107
Recall = 0.9612, SE=0.0081
f1 = 0.9518, SE=0.0078
Accuracy = 0.9556, SE=0.0067
```

# Making hard decisions

- logistic regression gives test score  $z = b + \mathbf{x}^T \mathbf{w}$ , which yields **confidence**  $\Pr\{y=1 \mid z\} = \frac{1}{1+e^{-z}}$
- Can convert to a **hard decision** by **thresholding**:
  - Set  $\hat{y} = 1$  when  $\Pr\{y=1 \mid z\} > t$  for **threshold**  $t$
- $t = 0.5$  minimizes the error rate,  $\Pr\{\hat{y} \neq y\}$ 
  - i.e., maximizes accuracy
- $t \in [0, 1]$  trades precision for recall
  - precision:  $\Pr\{y=1 \mid \hat{y}=1\}$
  - recall:  $\Pr\{\hat{y}=1 \mid y=1\}$
- $t \in [0, 1]$  trades sensitivity for specificity
  - sensitivity:  $\Pr\{\hat{y}=1 \mid y=1\}$
  - specificity:  $\Pr\{\hat{y}=0 \mid y=0\}$



# The ROC curve and the AUC

- The **receiver operating characteristic (ROC)** is the plot of TPR versus FPR

- **TPR**:  $\Pr\{\hat{y}=1 \mid y=1\}$

- **FPR**:  $\Pr\{\hat{y}=1 \mid y=0\}$

where each depends on the decision threshold  $t$

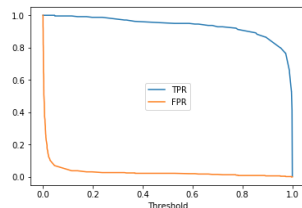
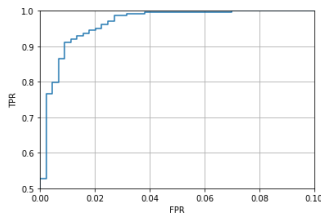
- Thus,  $t$  controls the location on the ROC curve

- The **area under the curve (AUC)** is a threshold-independent performance metric

```
from sklearn import metrics
yprob = logreg.predict_proba(Xs)
fpr, tpr, thresholds = metrics.roc_curve(y, yprob[:,1])
```

```
auc=metrics.roc_auc_score(y, yprob[:,1])
print("AUC=%f" % auc)
```

AUC=0.996344



# Performance metrics for multiclass classification

- In multiclass classification, there are many possible error types

- If the rows of the confusion matrix are normalized to sum-to-one ...

- $(k, l)$ th entry becomes

$$\Pr\{\hat{y}=l \mid y=k\}$$

- diagonal terms show per-class accuracy,  $\Pr\{\hat{y}=k \mid y=k\}$

contingency table or confusion matrix

	$\hat{y}=1$	$\hat{y}=2$	$\dots$	$\hat{y}=K$
$y=1$	10	3	$\dots$	4
$y=2$	1	14	$\dots$	3
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$y=K$	4	2	$\dots$	7

- The overall **accuracy** can be computed by averaging the per-class accuracies:

$$\Pr\{\hat{y} = y\} = \sum_{k=1}^K \Pr\{\hat{y}=k \mid y=k\} \Pr\{y=k\}$$

- Using **micro-averaging of  $k$ -vs-rest metrics**, can **generalize precision & recall**:

$$\text{precision} = \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \text{FP}_k} \quad \& \quad \text{recall} = \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \text{FN}_k}$$

and **generalize ROC** (i.e., the plot of TPR vs FPR) to multiclass data:

$$\text{TPR} = \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \text{FN}_k} \quad \& \quad \text{FPR} = \frac{\sum_k \text{FP}_k}{\sum_k \text{FP}_k + \text{TN}_k}$$

# Imbalanced class labels

- **Class imbalance** refers to the case where some classes have many fewer samples than others, i.e.,  $\Pr\{y=k\} \ll \frac{1}{K}$  for one or more  $k$
- In this case, accuracy is not a useful performance metric
  - It's possible to obtain high total accuracy  $\Pr\{\hat{y} = y\}$  even when per-class accuracy  $\Pr\{\hat{y} = k | y = k\}$  is very low for some  $k$ ! This can be seen from
 
$$\Pr\{\hat{y} = y\} = \sum_{k=1}^K \Pr\{\hat{y} = k | y = k\} \Pr\{y = k\}$$
  - Instead, use **balanced accuracy**  $\sum_{k=1}^K \Pr\{\hat{y} = k | y = k\}$ , or **AUC-PRC**
- Similarly, when training via cross-entropy loss, the minority classes are under-represented because they account for relatively few samples  $i$  in  $\sum_i \text{loss}_i$ 
  - Remedy: train using `class_weight='balanced'` (read more [here](#))
- Finally, when doing K-fold CV, use the **stratified** variant, which ensures that all folds have a similar class balance

# Learning objectives

- Understand classification problems in machine learning:
  - Identify features, labels; binary vs multiclass; linear vs nonlinear
  - visualize scatterplots and decision regions
- For binary classification problems, understand ...
  - linear classifiers, separating hyperplanes, linearly separable data
  - effect of feature transformation
  - why LS linear regression doesn't work well
  - logistic regression: logistic function, cross-entropy loss, ML fitting, regularization
  - common error metrics: accuracy, precision, recall, F1
  - the effect of the decision threshold, ROC, AUC
- For multiclass classification problems, understand ...
  - solutions that use multiple binary classifiers
  - multinomial logistic regression: softmax function, cross-entropy loss, ML fitting
- How to implement and assess classification using **sklearn**