

Unit 5

Linear Classification & Logistic Regression

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 4300: Introduction to Machine Learning, Sp20

Learning objectives

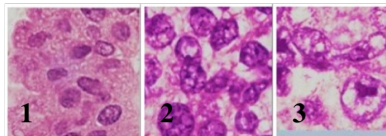
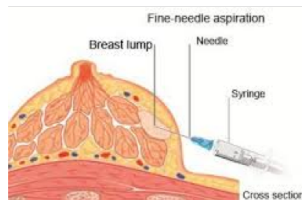
- Understand classification problems in machine learning:
 - Identify features, labels; binary vs multiclass; linear vs nonlinear
 - visualize scatterplots and decision regions
- For binary classification problems, understand ...
 - linear classifiers, separating hyperplanes, margin
 - why linear regression doesn't work
 - logistic regression: logistic function, cross-entropy loss, ML fitting, regularization
 - feature transformations
 - 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, 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
 - Tries to classify as **benign** or **malignant**
 - If malignant, also provides grading
- Would like to improve accuracy
- 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:

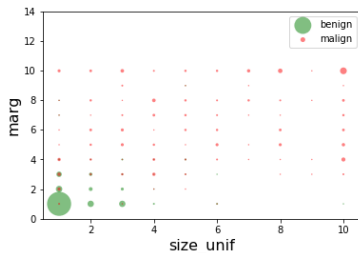
<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.names>

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
- Plot target count vs. features using variable-radius dots



- 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**
- Using $\{-1, 1\}$ leads to cleaner notation later

- Many applications:

- Face detection: is a face present here or not?
- Are these cells cancerous or not?

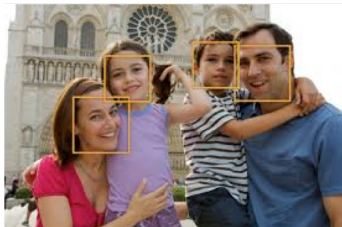
- Mathematically, want to design a **classifier**

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

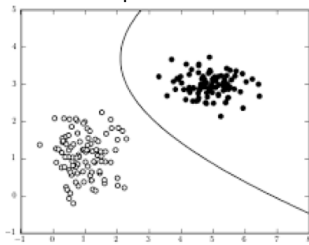
such that $\hat{y} = y$ with high probability

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

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



Example with $d = 2$:



Binary linear classification

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

- 1) define the “**score**” or “**discriminant**”, $z = b + \sum_{j=1}^d x_j w_j$, which is linear in the parameters b and w_j
- 2) given \mathbf{x} , threshold the score to obtain $\hat{y} = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$
- 3) learn weights $\mathbf{w} \triangleq [w_1, \dots, w_d]^T$ and intercept b from the training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

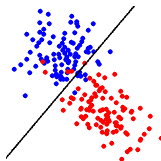
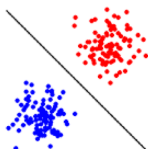
- Note that, at the **decision boundary**, we have

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

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 performs well when the data $\{(x_i, y_i)\}_{i=1}^n$ is “linearly separable”

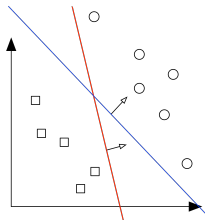
- Linearly separable means that there exists a hyperplane

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

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

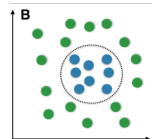
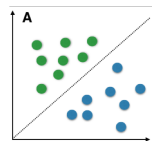
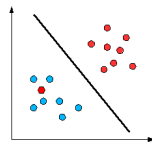
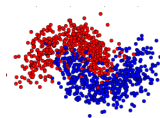
- On one side of hyperplane lies all x_i for which $y_i = 1$, while on other other side lies all x_i for which $y_i = -1$
- Note: When such a separating hyperplane exists, it is usually not unique.

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., shows which features matter)
 - 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 $\sqrt{x_1^2 + x_2^2}$
 - It can be used as a building block (e.g., for neural networks, decision trees, etc.)



Linear classification vs. linear regression

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

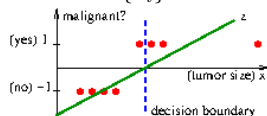
$$\hat{y} = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}, \quad \text{where } z = [1 \ x]^T \begin{bmatrix} b_{\text{ls}} \\ w_{\text{ls}} \end{bmatrix}$$

- **No, this may not work well!**
 - Consider simple case where $d = 1$
 - When $\{x_i\}$ is “balanced” (see figure), works okay
 - But when $\{x_i\}$ is “imbalanced,” the least-squares regression line gets pulled to one side, and the threshold at $z = 0$ will make errors
- What’s the problem?
 - **RSS is not the right loss function** for classification! (More details soon)

balanced $\{x_i\}$:



imbalanced $\{x_i\}$:



Linear regression fails on breast-cancer classification!

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

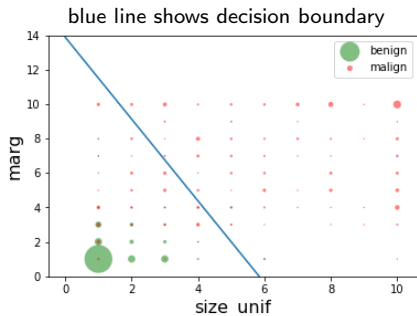
- Again, use decision boundary

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

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) via linear regression 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

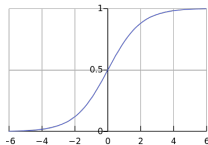
Binary logistic regression

- **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 is a bad idea!
- Idea: Given the score z , model the true label $y \in \pm 1$ as a **random variable**
- The most popular version of this uses

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

Note $\Pr\{y=0 \mid z\} + \Pr\{y=1 \mid z\} = 1 \ \forall z$, as required for a valid pmf.

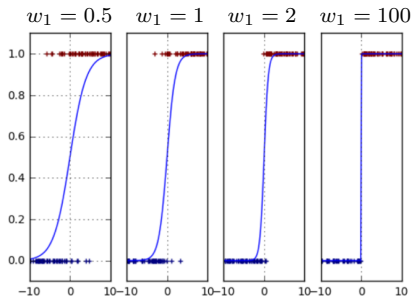
- The larger that z is, the more likely that $y = 1$
- When $z = 0$, it's equally likely that $y = 1$ or $y = -1$
- $\frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$ is known as the **“logistic function”** or **“sigmoid function”**



Understanding the logistic model

- Previously considered $\Pr\{y=1 \mid z\}$.
What about $\Pr\{y=1 \mid x\}$?
- Consider the simple case of a single scalar feature x and no intercept b :

$$\Pr\{y=1 \mid x\} = \frac{1}{1 + e^{-z}} \text{ for } z = w_1 x$$
 - As the weight w_1 increases, the transition becomes sharper
 - Equivalently, as w_1 increases, there is less randomness in y



curve: $\Pr\{y=1 \mid x\}$ versus x
 markers: random $\{(y_i, x_i)\}_{i=1}^n$

- Now add the intercept term b , giving

$$z = b + w_1 x = w_1 \left(\frac{b}{w_1} + x \right)$$

So the transition occurs at $x = -\frac{b}{w_1}$

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** $p(\mathbf{y}|\mathbf{X}; b, \mathbf{w})$ with model parameters (b, \mathbf{w})

- 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

$$\begin{aligned} (b_{\text{ml}}, \mathbf{w}_{\text{ml}}) &\triangleq \arg \max_{b, \mathbf{w}} p(\mathbf{y}|\mathbf{X}; b, \mathbf{w}) \\ &= \arg \max_{b, \mathbf{w}} \ln p(\mathbf{y}|\mathbf{X}; b, \mathbf{w}) \\ &= \arg \min_{b, \mathbf{w}} \{ -\ln p(\mathbf{y}|\mathbf{X}; b, \mathbf{w}) \} \end{aligned}$$

- In fact, we used the ML approach for linear regression:

- There, the likelihood was $p(\mathbf{y}|\mathbf{X}; \beta) = \mathcal{N}(\mathbf{y}; \mathbf{A}\beta, \sigma_\epsilon^2 \mathbf{I})$ with $\mathbf{A} = [\mathbf{1} \ \mathbf{X}]$

- thus $-\ln p(\mathbf{y}|\mathbf{X}; \beta) = \frac{1}{2\sigma_\epsilon^2} \|\mathbf{y} - \mathbf{A}\beta\|^2 + \text{constant}$

- and so $\beta_{\text{ml}} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{A}\beta\|^2 = \arg \min_{\beta} \text{RSS}(\beta) = \beta_{\text{ls}}$

- For logistic regression, we use the ML approach with a different likelihood!

ML estimation for logistic regression

- Logistic regression assumes that y_i depends only on \mathbf{x}_i (not $\mathbf{x}_k|_{k \neq i}$), and that

$$\Pr\{y_i = 1 \mid \mathbf{x}_i; b, \mathbf{w}\} = \frac{e^{z_i}}{1 + e^{z_i}}, \quad \Pr\{y_i = -1 \mid \mathbf{x}_i; b, \mathbf{w}\} = \frac{1}{1 + e^{z_i}}, \quad z_i = b + \mathbf{x}_i^\top \mathbf{w}$$

- Thus we have

$$\begin{aligned}
 p(\mathbf{y} \mid \mathbf{X}; b, \mathbf{w}) &= \prod_{i=1}^n p(y_i \mid \mathbf{x}_i; b, \mathbf{w}) && \text{via independence assumption} \\
 -\ln p(\mathbf{y} \mid \mathbf{X}; b, \mathbf{w}) &= -\sum_{i=1}^n \ln p(y_i \mid \mathbf{x}_i; b, \mathbf{w}) && \text{since } \ln(ab) = \ln a + \ln b \\
 &= -\sum_{i=1}^n \frac{y_i + 1}{2} \ln \Pr\{y_i = 1 \mid \mathbf{x}_i; b, \mathbf{w}\} + \left(1 - \frac{y_i + 1}{2}\right) \ln \Pr\{y_i = -1 \mid \mathbf{x}_i; b, \mathbf{w}\} \\
 &\quad \text{since } \frac{y_i + 1}{2} = \begin{cases} 1 & y_i = 1 \\ 0 & y_i = -1 \end{cases} \\
 &= -\sum_{i=1}^n \frac{y_i + 1}{2} (z_i - \ln[1 + e^{z_i}]) + \left(1 - \frac{y_i + 1}{2}\right) (0 - \ln[1 + e^{z_i}]) \\
 &= -\sum_{i=1}^n \left(\frac{y_i + 1}{2} z_i - \frac{y_i + 1}{2} \ln[1 + e^{z_i}] - \ln[1 + e^{z_i}] + \frac{y_i + 1}{2} \ln[1 + e^{z_i}] \right) \\
 &= -\sum_{i=1}^n \left(\frac{y_i + 1}{2} z_i - \ln[1 + e^{z_i}] \right)
 \end{aligned}$$

- The ML estimates of (b, \mathbf{w}) are the values that minimize this expression

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 commonly occurring alternative

- 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 cross-entropy loss:

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

- Why? When the training data is linearly separable, the cross-entropy loss decreases as $|w_j|$ increases, and so $\|\mathbf{w}_{ml}\| = \infty$. L2 regularization keeps \mathbf{w}_{lr} finite
 - Could instead use **L1 regularization**, and thus perform feature selection:
- $$(b_{lr}, \mathbf{w}_{lr}) = \arg \min_{b, \mathbf{w}} \left\{ \sum_{i=1}^n (\ln[1 + e^{z_i}] - y_i z_i) + \alpha \|\mathbf{w}\|_1 \right\} \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$
- With L2 regularization, could use small (but positive) α to avoid bias
 - Or, with L1 or L2 regularization, could tune α using **K-fold cross-validation**

Logistic regression in sklearn

In `sklearn`, there is a nice `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 use *large* C to avoid regularization-induced bias

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

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

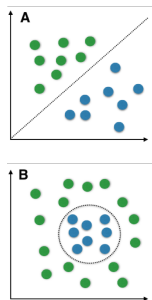
```
LogisticRegression(C=100000.0, class_weight=None, dual=False,
                    fit_intercept=True, intercept_scaling=1, max_iter=100,
                    multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
                    solver='liblinear', tol=0.0001, verbose=0, warm_start=False)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

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

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

Feature transformations

- Can we implement *curved* boundaries with *linear* classification methods like logistic regression?
- Yes! Through **feature transformation** ...
 - Example: given raw features (x_1, x_2) , we could transform to $\mathbf{x} \triangleq [x_1, x_2, x_1^2, x_2^2]^T$
 - Then the “linear” score $z = b + \mathbf{x}^T \mathbf{w}$ would be *quadratic* in the raw features (x_1, x_2)
- **One-hot coding** is another important feature transformation. We saw how it can be used for categorical features like $x_j \in \{\text{Ford, BMW, GM}\}$



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 $K \geq 2$ classes?
 - Binary classification is the special case $K = 2$

- 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**

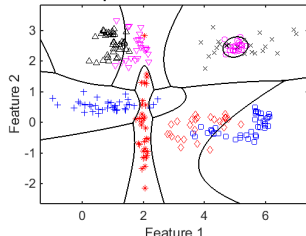
$$f(x) = \hat{y} \in \{1, \dots, K\}$$

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

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

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

example where $K = 10$:



Multiclass classification using binary classifiers

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

Two main approaches:

■ Voting on binary outcomes

- **1-vs-all:** For each class k , decide if x is **in- k vs not-in- k** , then choose the k with most votes
- **1-vs-1:** For each pair (k, l) , decide if x is **in- k vs in- l** , then choose k with most votes

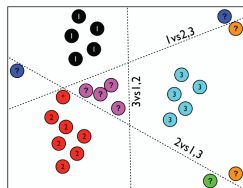
Problem: Tied votes lead to ambiguities! (see right)

■ Choosing the highest confidence value

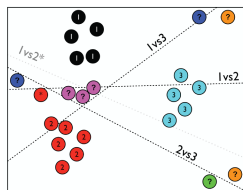
- **1-vs-all:** For each class k , compute confidence that x is **in- k vs not-in- k** , then choose the k with highest confidence

Ambiguities avoided by continuous-valued confidence!

We'll study one example (MLR) in more detail ...



(b) 1-vs-All



(a) 1-vs-1

Outline

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

Multinomial logistic regression

Suppose we combine **multiple binary classifiers** using the “1-vs-all confidence-value” approach (p. 25). With binary *linear* classifiers, we would do:

- For each class $k = 1, \dots, K$,
 - design b_k and $\mathbf{w}_k = [w_{k1}, \dots, w_{kd}]^T$ to classify **in- k or not-in- k**
 - compute a linear **score** from \mathbf{x} , i.e., $z_k = b_k + \sum_{j=1}^d x_j w_{kj}$

then choose the winner via: $\hat{y} = \arg \max_k \underbrace{\Pr\{y=k \mid \mathbf{z}\}}_{\text{confidence value}}$

- In **multinomial logistic regression**, we adopt the “**softmax**” likelihood function:

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

- The softmax has the property that, if $z_{k_{\max}} \gg z_k$ for all $k \neq k_{\max}$, then

$$\frac{e^{z_k}}{\sum_{l=1}^K e^{z_l}} \approx \begin{cases} 1 & \text{if } k = k_{\max} \\ 0 & \text{if } k \neq k_{\max} \end{cases} \quad \text{so it's a “soft” version of “}\delta_{k-k_{\max}}\text{”}$$

One-hot label-coding for MLR

- We design $\mathbf{b} = [b_1, \dots, b_K]^\top$ & $\mathbf{W} \triangleq [\mathbf{w}_1, \dots, \mathbf{w}_K]$ using **ML estimation**, i.e.,

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

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

$$y_{ik} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{if } y_i \neq k \end{cases} \quad \text{for all } i = 1, \dots, n$$

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

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

Multinomial cross-entropy loss

- Combining the results from the previous page, we get

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

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

- Could also add $\alpha \|\mathbf{W}\|_F^2$ (L2 regularization) or $\alpha \|\mathbf{W}\|_1$ (L1 regularization), and tune α via cross-validation
- With or without regularization, there is no closed-form expression for optimal (\mathbf{b}, \mathbf{W}) , but the optimization problems are convex and can be solved numerically
- For this, **sklearn** provides an excellent **LogisticRegression** implementation

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

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

- The implications of these errors can be very different!

- e.g., consider breast cancer diagnosis

- Common machine-learning performance metrics:

- precision**: $\Pr\{y=1 \mid \hat{y}=1\} = \frac{TP}{TP+FP}$ given a positive test, how often is the patient cancerous?

- recall**: $\Pr\{\hat{y}=1 \mid y=1\} = \frac{TP}{TP+FN}$ given that the patient has cancer, how often is it detected?

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

- accuracy**: $\Pr\{\hat{y}=y\} = \frac{TP+TN}{TP+FN+TN+FP}$ how often is the test correct?

- Common metrics in medicine:

- sensitivity**: $\Pr\{\hat{y}=1 \mid y=1\} = \frac{TP}{TP+FN}$ given that the patient has cancer, how often is it detected?

- specificity**: $\Pr\{\hat{y}=0 \mid y=0\} = \frac{TN}{TN+FP}$ given a healthy patient, how often is diagnosis correct?

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.metrics import precision_recall_fscore_support

nfold = 10
kf = KFold(n_splits=nfold, shuffle=True)
prec = []
rec = []
f1 = []
acc = []
for train, test in kf.split(Xs):
    # Get training and test data
    Xtr = Xs[train,:]
    ytr = y[train]
    Xts = Xs[test,:]
    yts = y[test]

    # Fit a model
    logreg.fit(Xtr, ytr)
    yhat = logreg.predict(Xts)

    # Measure performance
    preci, reci, f1i, _ = precision_recall_fscore_support(yts, yhat, average='binary')
    prec.append(preci)
    rec.append(reci)
    f1.append(f1i)
    acci = np.mean(yhat == yts)
    acc.append(acci)

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

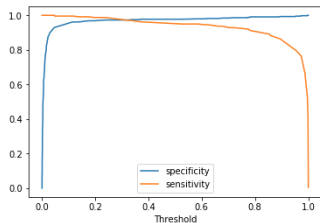
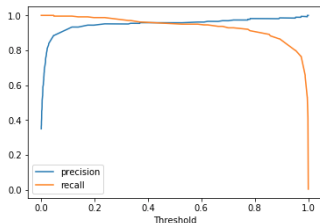
# Compute the standard errors
prec_se = np.std(prec)/np.sqrt(nfold-1)
rec_se = np.std(rec)/np.sqrt(nfold-1)
f1_se = np.std(f1)/np.sqrt(nfold-1)
acc_se = np.std(acc)/np.sqrt(nfold-1)

Precision = 0.9554, SE=0.0095
Recall = 0.9513, SE=0.0095
f1 = 0.9527, SE=0.0051
Accuracy = 0.9678, SE=0.0037

```


Making hard decisions

- Logistic regression outputs a **confidence value**, $\Pr\{y=1 \mid \mathbf{x}\} \in [0, 1]$
- Can convert to a **hard decision** by **thresholding**:
 - Set $\hat{y} = 1$ when $\Pr\{y=1 \mid \mathbf{x}\} > t$ for **threshold** t
- $t = 1/2$ 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: $\Pr\{\hat{y}=0 \mid y=0\}$



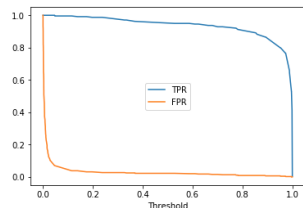
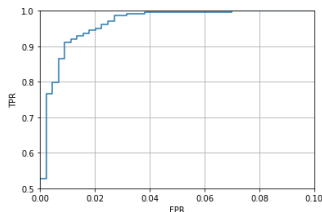
The ROC curve and the AUC

- The **receiver operating characteristic (ROC)** is the plot of TPR versus FPR
 - **FPR**: $\Pr\{\hat{y}=1 \mid y=0\}$
 - **TPR**: $\Pr\{\hat{y}=1 \mid y=1\}$
- The threshold t controls location on 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 columns of confusion matrix are normalized to sum-to-one ...

- (k, l) th entry becomes

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

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

- The overall accuracy can be computed as

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

contingency table or confusion matrix

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

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, margin
 - why linear regression doesn't work
 - logistic regression: logistic function, cross-entropy loss, ML fitting, regularization
 - feature transformations
 - 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, ML fitting
- How to implement and assess classification using **sklearn**