

Udacity: Intro to ML with PyTorch

Contents

1	Supervised Learning	1
1.1	Linear Regression (doc)	1
1.2	Perceptron Algorithm	3
1.3	Decision Trees (doc)	3
1.4	Naive Bayes (doc)	5
1.5	Support Vector Machines (SVMs) (doc)	6
1.6	Ensemble Methods (doc)	7
1.7	Model Evaluation Metrics	8
1.7.1	Classification Metrics (formulas)	8
1.7.2	Regression Metrics	9
1.8	Training and Tuning	10
2	Deep Learning	12
2.1	Neural Networks	12
2.2	Implementing Gradient Descent (NumPy)	13
2.3	Training Neural Networks	14
2.4	Deep Learning with PyTorch	15
2.4.1	Building Neural Networks	15
2.4.2	Training Neural Networks	17
2.4.3	Validation and Dropout	18
2.4.4	Saving and Loading Models	19
2.4.5	Loading Image Data	19
2.4.6	Transfer Learning	20
3	Unsupervised Learning	23
3.1	Clustering (KMeans)	23
3.2	Hierarchical and Density Based Clustering	23
3.2.1	Hierarchical Clustering	23
3.2.2	Density-Based Clustering	24
3.3	Gaussian Mixture Models and Cluster Validation	25
3.3.1	Cluster Analysis Process	25
3.4	Dimensionality Reduction and PCA	26
3.5	Random Projections and ICA	28

1 Supervised Learning

1.1 Linear Regression (doc)

Absolute Trick: A point at (p, q) with a line $y = w_1x + w_2$, we want to change the line to slowly move it towards the point. We want to do this by a small value α (learning rate), so our line become $\hat{y} = (w_1 + p\alpha)x + (w_2 + \alpha)$.

Square Trick: Similar to the absolute trick, but we consider the offset of the point from the y-axis (p) and from the line ($q - q'$) where q' is the *estimated* y point and q is the actual y point. So our line becomes $\hat{y} = (w_1 + p(q - q')\alpha)x + (w_2 + (q - q')\alpha)$

Mean Absolute Error: the average of the sum of the distance between the estimate line value \hat{y} and actual y value, giving us the equation: $Error = \frac{1}{m} \sum_{i=1}^m |y - \hat{y}|$

Mean Squared Error: Similar to MAE, but we square the value instead which gives us a square near the value $y - \hat{y}$. This gives us the equation: $Error = \frac{1}{2m} \sum_{i=1}^m (y - \hat{y})^2$

Gradient Step: $w_i \rightarrow w_i - \alpha \frac{\partial}{\partial w_i} Error$ with $\frac{\partial}{\partial w_1} Error = -(y - \hat{y})x$ and $\frac{\partial}{\partial w_2} Error = -(y - \hat{y})$

```
1 X : array of predictor features
2 y : array of outcome values
3 W : predictor feature coefficients
4 b : regression function intercept
5 learn_rate : learning rate
6
7 yhat = np.matmul(X,W) + b # predicted values
8 error = y-yhat
9 W_new = W + learn_rate*np.matmul(error,X) # updated slope
10 b_new = b + learn_rate*error.sum() # updated y-intercept
```

Higher Dimensions: $\hat{y} = w_1x_1 + w_2x_2 + \dots + w_{n-1}x_{n-1} + w_n$

Polynomial Regression: Allows us to add degrees to our x to better find our line.

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import LinearRegression
5
6 train_data = pd.read_csv('data.csv')
7 X = np.asarray(train_data['Var_X']).reshape((20,1))
8 y = np.asarray(train_data['Var_Y']).reshape((20,1))
9
10 poly_feat = PolynomialFeatures(degree=4) # 4th degree polynomial
11 X_poly = poly_feat.fit_transform(X)
12
13 poly_model = LinearRegression(fit_intercept=False).fit(X_poly, y)
```

L1 Regularization: adds the absolute value of the coefficients (w_i) to the error. This is computationally inefficient, better when data is sparse, gives us feature selection (makes irrelevant columns to 0).

L2 Regularization: we add the square of the coefficients (w_i) to the error. This is computationally efficient, better for non-sparse data, doesn't have feature selection (treats all columns equally).

λ **Parameter:** determines how much impact the coefficients have on our total error. Choosing a large λ leads to simpler models more while a small λ leads to more complex models.

We can train a Linear Regression model with L1 regularization applied to it by using the [Lasso model](#).

```
1 import pandas as pd
2 from sklearn.linear_model import Lasso
3
4 train_data = pd.read_csv('data.csv', header=None)
5 X = train_data.iloc[:, :-1] # all rows, all columns up to last
6 y = train_data.iloc[:, -1] # all rows, last column
7
8 lasso_reg = Lasso().fit(X, y)
9
10 reg_coef = lasso_reg.coef_
11 print(reg_coef)
12 # [ 0.  2.35793224  2.00441646 -0.05511954 -3.92808318  0.]
```

Feature Scaling: a way of transforming your data into a common range of values (2 ways below).

Standardizing: turning the column into a standard normal variable by computing $\frac{x-\mu}{\sigma}$ and is interpreted as the number of standard deviations each value is away from the mean (most common type).

Normalizing: data is scaled between 0 and 1 by calculating $\frac{x-x.min}{x.max-x.min}$

Distance Based Metrics: If an algorithm uses distance based metrics to predict, choosing some sort of feature scaling is necessary (such as for SVMs or k-nn).

When using Regularization: we need all columns to have equal ranges so they are treated similar when calculating the penalty. We need features with small ranges and large ranges to have similar coefficients.

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.linear_model import Lasso
4
5 train_data = pd.read_csv('data.csv', header=None)
6 X = train_data.iloc[:, :-1]
7 y = train_data.iloc[:, -1]
8
9 scaler = StandardScaler()
10 X_scaled = scaler.fit_transform(X)
11
12
13 lasso_reg = Lasso().fit(X_scaled, y)
14
15 reg_coef = lasso_reg.coef_
16 print(reg_coef)
17 # [0.  3.90753617  9.02575748 -0.  -11.78303187  0.45340137]
```

Notice how after we scaled the features, the L1 regularization set the 1st and 4th columns coefficients to zero. Where when we did not standardize in the previous code, the L1 regularization set the 1st and 6th columns coefficients to zero.

1.2 Perceptron Algorithm

Linear Boundaries: we use a boundary line to determine how to classify the data, denoted by the equation $w_1x_1 + w_2x_2 + b = 0$ which we rewrite as $Wx + b = 0$ where $W = (w_1, w_2)$ and $x = (x_1, x_2)$. Note that we just take the dot product of $W \cdot x \implies w_1x_1 + w_2x_2$. We then classify the output of the boundary line by:

$$\hat{y} = \begin{cases} 1 & \text{if } Wx + b \geq 0 \\ 0 & \text{if } Wx + b < 0 \end{cases}$$

Higher Dimensions: if we are working with n columns, our boundary line becomes a $n - 1$ dimensional hyperplane with the equation $w_1x_1 + w_2x_2 + \dots + w_nx_n + b = 0 \implies Wx + b = 0$. The prediction works the same way as the linear boundary. Note that weights (W) is a $(1 \times n)$ row vector, input (x) is a $(n \times 1)$ column vector, and bias (b) is a (1×1) vector.

Basic Perceptrons: We have our input nodes (x_1, \dots, x_n and can also include our bias b) which are then multiplied by the corresponding weights (w_1, \dots, w_n), which are then inputted into a linear function node ($Wx + b$), and finally that input is sent to a step function node (determines if \hat{y} is 1 or 0).

Updating the line: In a 2D example, with point (p, q) , we can update the boundary line with the equation $(w_1 \pm \alpha p)x_1 + (w_2 \pm \alpha q)x_2 + (b \pm \alpha * 1)$. Note that \pm depends on which way we move the line.

```
1 def perceptronStep(X, y, W, b, learn_rate = 0.01):
2     for i in range(len(X)): # go through all points
3         yhat = prediction(X[i], W, b) # get predicted values
4         if y[i] - yhat == 1: # if y=1 but yhat=0 (move line down)
5             W[0] += learn_rate*X[i][0]
6             W[1] += learn_rate*X[i][1]
7             b += learn_rate
8         if y[i] - yhat == -1: # if y=0 but yhat=1 (move line up)
9             W[0] -= learn_rate*X[i][0]
10            W[1] -= learn_rate*X[i][1]
11            b -= learn_rate
12    return W, b
```

1.3 Decision Trees (doc)

(Multiclass) Entropy: $-\sum_{i=1}^n p_i \log_2(p_i)$ where p_i is the probability for each class. Note that points with the same color will have zero entropy, so this is the ideal outcome for classification.

Information Gain: Tells us how well our data is being classified (in a range between 0 and 1). The lower the information gain, the worse our model is doing as separating the data. The higher the information gain, the better our model is performing at classification. We have the formula:

$$IG = Entropy(Parent) - \left[\frac{m}{m+n} Entropy(Child_1) + \frac{n}{m+n} Entropy(Child_2) \right]$$

Important Hyperparameters for Decision Trees:

Maximum depth: the largest possible length between the root to a leaf, a tree of maximum length k can have at most 2^k leaves.

Minimum samples to split: numeric value (m) that says if a node doesn't have at least m samples it can't split.

Minimum samples per leaf: decides how many samples must be in each node. If it is an integer, then it is a sample size. If it is a float, it is a percentage of samples needed in a leaf.

WARNING: Large depth often causes overfitting. Small depth can result in a very simple model, which may cause underfitting. Small minimum samples per split may result in overfitting. Large minimum samples may result in the tree not having enough flexibility to get built, and may result in underfitting.

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.metrics import accuracy_score
6 import random
7 random.seed(42)
8
9 in_file = 'titanic_data.csv'
10 full_data = pd.read_csv(in_file)
11
12 outcomes = full_data['Survived'] # Labels
13 features_raw = full_data.drop('Survived', axis = 1) # Remove label from features df
14
15 features_no_names = features_raw.drop(['Name'], axis=1) # Remove names
16
17 features = pd.get_dummies(features_no_names) # One-hot encoding
18 features = features.fillna(0.0)
19
20 len(features.columns) # 839
21 # Note that we now have 839 columns instead of 11 since each ticket is a column
22 # from the one-hot encoding (each ticket_number is now a column with all 0's
23 # expect for the person/row with the ticket, which will be a 1).
24
25 X_train, X_test, y_train, y_test = train_test_split(features, outcomes,
26                                                    test_size=0.2, random_state=42)
27
28 model = DecisionTreeClassifier().fit(X_train, y_train)
29
30 # Making predictions
31 y_train_pred = model.predict(X_train)
32 y_test_pred = model.predict(X_test)
33 train_accuracy = accuracy_score(y_train, y_train_pred) # 1.0
34 test_accuracy = accuracy_score(y_test, y_test_pred) # 0.8156
35 # High training and lower test accuracy means we might be overfitting.
36
37 # Tune the hyperparameters for new model to improve test score
38 model_tune = DecisionTreeClassifier(max_depth=20, min_samples_leaf=6)
39 model_tune.fit(X_train, y_train)
40 tune_y_train_pred = model_tune.predict(X_train)
41 tune_y_test_pred = model_tune.predict(X_test)
42 tune_train_accuracy = accuracy_score(y_train, tune_y_train_pred)
43 tune_test_accuracy = accuracy_score(y_test, tune_y_test_pred)
44
45 print('Train Acc: {}, Test Acc: {}'.format(tune_train_accuracy, tune_test_accuracy))
46 # Train Acc: 0.8820224719101124, Test Acc: 0.8603351955307262

```

1.4 Naive Bayes (doc)

Naive Bayes: a supervised machine learning algorithm that can be trained to classify data into multi-class categories. The probabilistic model computes the conditional probabilities of the input features and assigns the probability distributions to each of possible classes.

Bayes Theorem: We have prior probabilities for a given problem, and after we know that event B occurred, we can find the posterior probability with the formula $P(A_j|B) = \frac{P(B|A_j)P(A_j)}{\sum P(B|A_i)P(A_i)}$

Naive Assumption: $P(A \cap B) = P(A)P(B)$ meaning that the two events are independent. We assume this to be true even if it is not the case.

Conditional Probability (Proportional): $P(A|B_1, B_2, \dots, B_n) \propto P(B_1, B_2, \dots, B_n|A)P(A)$

Algorithm: For an event A , we have $P(A|B_1, B_2, \dots, B_n) \propto \frac{P(B_1|A)P(B_2|A)\dots P(B_n|A)P(A)}{P(B_1, B_2, \dots, B_n)}$

Bag of Words (BoW): a collection of text data where we count the frequency of the words.

Accuracy: the ratio of the number of correct predictions to the total number of predictions.

Precision: a ratio of true positives to all positives (true+false positive).

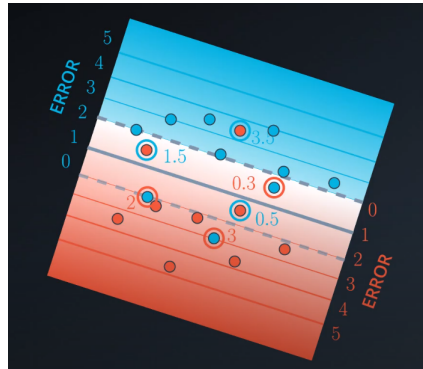
Recall (Sensitivity): a ratio of true positives to all the words (true positives + false negatives).

F1 Score: is the weighted average of the precision and recall scores (between 0 and 1, higher better).

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Read in data from table, add column names, specify separator
8 df = pd.read_table('smsspamcollection/SMSSpamCollection', sep='\t',
9                   names=['label', 'sms_message'])
10
11 # Convert labels to numerical encodings
12 df['label'] = df.label.map({'ham':0, 'spam':1})
13
14 X_train, X_test, y_train, y_test = train_test_split(df['sms_message'], df['label'],
15                                                    random_state=1)
16
17 count_vector = CountVectorizer()
18 # Note that CountVectorizer automatically converts to lowercase and ignores
19 # all punctuation. We can add stop_words=English if we want to ignore
20 # common english words (a, an, and) in our matrix.
21
22 # Fit the training data and return the frequency matrix
23 training_data = count_vector.fit_transform(X_train)
24
25 # Transform testing data and return the matrix (don't fit testing data)
26 testing_data = count_vector.transform(X_test)
27
28 naive_bayes = MultinomialNB()
29 naive_bayes.fit(training_data, y_train)
30 predictions = naive_bayes.predict(testing_data)
31
32 print('Accuracy score: ', format(accuracy_score(y_test, predictions))) # 0.98851
33 print('Precision score: ', format(precision_score(y_test, predictions))) # 0.972067
34 print('Recall score: ', format(recall_score(y_test, predictions))) # 0.94054
35 print('F1 score: ', format(f1_score(y_test, predictions))) # 0.95604
```

1.5 Support Vector Machines (SVMs) (doc)

Classification Error: we have a line separating our data points into groups, $Wx + b = 0$, and we want to add two more lines called the *margin*. The line above has equation $Wx + b = 1$ and the line below has equation $Wx + b = -1$. We then assign points error based on where they are in respect to the margin:



Classification Error: we take the absolute value of all misclassified points and add them together.

Margin: we want a small error if the margin is large, and a large error if the margin is small. We can determine the size of the margin by $\frac{2}{|W|}$ (2 divided by the norm of W).

Margin Error: $|W|^2$ (note this is the same value given by L2 regularization).

Error Function: we want $\text{Error} = \text{Classification Error} + \text{Margin Error}$, which we will minimize using gradient descent.

The C Parameter: $\text{Error} = C * \text{Classification Error} + \text{Margin Error}$. It is a constant we use to determine the influence of the classification error. Large C means focus on classifying points (smaller margin), small C means focus on margin error (makes more classification errors).

Linear Kernel: we separate our data using a linear line (doesn't work well when data is complicated).

Polynomial Kernel: we separate our data using circles, hyperbolas, parabolas, and many more boundaries by adding more dimensions to our data to classify it. The degree of the polynomial determines what kind of boundaries we can create.

ex: (degree 2) $\mathbb{E}^2 \rightarrow \mathbb{E}^5 : (x, y) \rightarrow (x, y, x^2, xy, y^2)$. We map our points (x, y) to a higher dimensions which allows us to create degree 2 functions. We then find a boundary hyperplane and map this back to a 2nd degree polynomial boundary (think of a paraboloid in 3D mapping back to a circle in 2D).

RBK Kernel: We build mountain ranges (radial basis functions) on top of our point to find boundaries to group the data. The plane that intersects the paraboloids is the projection that determine the boundaries. The γ hyperparameter can be tuned to determine the width of the curve, with a small γ giving us a wide curve and a large γ giving us a narrow curve.

```
1 from sklearn.svm import SVC
2 from sklearn.metrics import accuracy_score
3
4 data = np.asarray(pd.read_csv('data.csv', header=None))
5 X = data[:,0:2] # features
6 y = data[:,2] # labels
7
8 model = SVC(kernel='rbf', gamma=27).fit(X,y)
9
10 y_pred = model.predict(X)
11 acc = accuracy_score(y, y_pred)
12
```

1.6 Ensemble Methods (doc)

Bias: when a model has high bias, this means that it doesn't do a good job of bending to the data. One example is linear regression, which fits a straight line no matter how the data is shaped (but has a low variance).

Variance: when a model has high variance, this means that it changes drastically to meet the needs of every point in our dataset. One example are decision trees, which will attempt to split every point into its own branch if possible (but has low bias).

High Bias, Low Variance: models tend to underfit data, as they are not flexible (linear models).

High Variance, Low Bias: models tend to overfit data, as they are too flexible (decision trees).

Randomness: introduction of randomness to high variance algorithms to combat the tendency of these algorithms to overfit the data.

Bootstrap: sampling the data with replacement and fitting your algorithm to the sampled data.

Boosting: Random sampling with replacement over weighted data.

Random Forests: take a subset of the columns and build a decision tree from them, then repeat this process multiple times. We then use each tree to predict, and pick the most occurring outcome.

Bagging: take subsets of the data, create one-node decision trees to find a boundary for each subset, then take the overlap of the boundaries to determine the final boundary line (bootstrapping).

AdaBoost: we assign weights to each points, and we want to minimize the sum of the weights on the incorrectly classified points. We then increase the weights (by $\frac{\sum \text{correct weights}}{\sum \text{incorrect weights}}$) of the misclassified points and create a new boundary line. We continue this process for multiple models.

Once we have our models, we assign weight to the models to determine the influence they have in our final boundary, which can be calculated by $weight = \ln(\frac{\# \text{correct}}{\# \text{incorrect}})$

```
1  # NOTE: we are using the testing and training data from the Naive Bayes
2  # section, so refer to that to see data preprocessing steps
3
4  from sklearn.ensemble import AdaBoostClassifier
5  from sklearn.ensemble import BaggingClassifier
6  from sklearn.ensemble import RandomForestClassifier
7  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
8
9  bag_model = BaggingClassifier(n_estimators=200).fit(training_data, y_train)
10 rf_model = RandomForestClassifier(n_estimators=200).fit(training_data, y_train)
11 ada_model = AdaBoostClassifier(n_estimators=300, learning_rate=0.2)
12 ada_model.fit(training_data, y_train)
13
14 bag_ypred = bag_model.predict(testing_data)
15 rf_ypred = rf_model.predict(testing_data)
16 ada_ypred = ada_model.predict(testing_data)
17
18 # General Format for printing evaluation metrics (scoring model)
19 print('Accuracy score: ', format(accuracy_score(y_true, preds)))
20 print('Precision score: ', format(precision_score(y_true, preds)))
21 print('Recall score: ', format(recall_score(y_true, preds)))
22 print('F1 score: ', format(f1_score(y_true, preds)))
23
```

NOTE: the *base_estimator* parameter for AdaBoost and Bagging allows you to choose the weak learner model to be used in the classification process.

1.7 Model Evaluation Metrics

1.7.1 Classification Metrics ([formulas](#))

Confusion Matrix: a table that tells us the True Positive, False Negatives, False Positive, and True Negatives (left to right, top to bottom).

Error Types: False Positives are classified as *Type 1 Error*, False Negative are *Type 2 Error*.

Accuracy: used to compare models, it tells us the proportion of observations we correctly labeled (don't use if classes are imbalanced).

Precision: focuses on *predicted* positive values, helps determine if you are doing a good job of predicting the positive values, as compared to predicting negative values as positive. Avoiding false positives.

Recall: focuses on the *actual* positive values, determines if you are doing a good job of predicting the positive values without regard of how you are doing on the actual negative values. Avoiding false negatives.

F1 Score: the harmonic mean (less than regular mean), found by $2 * \frac{Precision * Recall}{Precision + Recall}$ (care equally about both the positive and negative cases).

F_β Score: If we want our model to focus more on precision than recall (such as spam filtering), we can use a smaller β . If we want our model to focus more on recall than precision (such as medical diagnosis), we can use a larger β . The formula is $F_\beta = (1 + \beta^2) \frac{precision * recall}{\beta^2 * precision + recall}$

ROC Curve & AUC: By finding different thresholds for our classification metrics, we can measure the area under the ROC curve. When the AUC is closer to 1, the better your model is performing.

```
1  # Again, we are using the data from the Naive Bayes section, so refer
2  # back to that to see the data preprocessing steps
3
4  from sklearn.naive_bayes import MultinomialNB
5  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
6  fbeta_score
7  from sklearn.ensemble import BaggingClassifier, RandomForestClassifier,
8  from sklearn.svm import SVC
9
10 naive_bayes = MultinomialNB().fit(training_data, y_train)
11 bag_mod = BaggingClassifier(n_estimators=200).fit(training_data, y_train)
12 rf_mod = RandomForestClassifier(n_estimators=200).fit(training_data, y_train)
13 svm_mod = SVC().fit(training_data, y_train)
14
15 nb_preds = naive_bayes.predict(testing_data)
16 bag_preds = bag_mod.predict(testing_data)
17 rf_preds = rf_mod.predict(testing_data)
18 svm_preds = svm_mod.predict(testing_data)
19
20 def metric_scores(actual, preds, model, model_name):
21     print('Accuracy for ' + model_name + ' : {}'.format(accuracy_score(actual, preds)))
22     print('Precision ' + model_name + ' : {}'.format(precision_score(actual, preds)))
23     print('Recall for ' + model_name + ' : {}'.format(recall_score(actual, preds)))
24     print('F1 for ' + model_name + ' : {}'.format(f1_score(actual, preds)))
25     print('\n')
26     return None
27
28 metric_scores(y_test, nb_preds, naive_bayes, 'Naive Bayes')
29 metric_scores(y_test, bag_preds, bag_mod, 'Bagging')
30 metric_scores(y_test, rf_preds, rf_mod, 'Random Forest')
31 metric_scores(y_test, svm_preds, svm_mod, 'SVM')
32 # Running this, Naive Bayes performed best for all expect precision (rf_model)
33 fbeta_score(y_test, nb_preds, beta=1) # same output as F1 score since beta=1
```

```

1  from itertools import cycle
2  from sklearn.metrics import roc_curve, auc, roc_auc_score
3  from scipy import interp
4
5  def build_roc_auc(model, X_train, X_test, y_train, y_test):
6      y_preds = model.fit(X_train, y_train).predict_proba(X_test)
7      # Compute ROC curve and ROC area for each class
8      fpr = dict()
9      tpr = dict()
10     roc_auc = dict()
11     for i in range(len(y_test)):
12         fpr[i], tpr[i], _ = roc_curve(y_test, y_preds[:, i])
13         roc_auc[i] = auc(fpr[i], tpr[i])
14
15     # Compute micro-average ROC curve and ROC area
16     fpr["micro"], tpr["micro"], _ = roc_curve(y_test.ravel(), y_preds[:, 1].ravel())
17     roc_auc["micro"] = auc(fpr["micro"], tpr["micro"])
18
19     plt.plot(fpr[2], tpr[2], color='darkorange', lw=2,
20             label='ROC curve (area = %0.2f)' % roc_auc[2])
21     plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
22     plt.xlim([0.0, 1.0])
23     plt.ylim([0.0, 1.05])
24     plt.xlabel('False Positive Rate')
25     plt.ylabel('True Positive Rate')
26     plt.title('Receiver operating characteristic example')
27     plt.show();
28
29     return roc_auc_score(y_test, np.round(y_preds[:, 1]))

```

1.7.2 Regression Metrics

Mean Absolute Error (MAE): a useful metric to optimize on when the value you are trying to predict follows a skewed distribution (not as much influence from outliers as MSE). The optimal value for this technique is the median value.

Mean Squared Error (MSE): most used metric for optimization (impacted by outliers and skew), it is differentiable and can be used better for gradient descent. The optimal value is the mean.

R2 Score: interpreted as the ‘amount of variability’ captured by a model. Minimizing MSE will maximize R2, which is good since we want this score to be as close to 1 as possible.

```

1  # Import all regression models
2  from sklearn.linear_model import LinearRegression
3  from sklearn.ensemble import RandomForestRegressor, AdaBoostRegressor
4  from sklearn.tree import DecisionTreeRegressor
5  from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
6
7  def reg_metrics(actual, preds, model, model_name):
8      print('R2 ' + model_name + ' : {}'.format(r2_score(actual, preds)))
9      print('MAE ' + model_name + ' : {}'.format(mean_absolute_error(actual, preds)))
10     print('MSE ' + model_name + ' : {}'.format(mean_squared_error(actual, preds)))
11     print('\n')
12     return None

```

Using the Boston housing dataset, if we were to fit and predict on all of the models above, we would see that the RandomForestRegressor performs the best on all evaluation metrics.

1.8 Training and Tuning

Underfitting: making our model too simple, leading to bad accuracy on the training and testing set (an error due to high bias). Does not generalize well to new data.

Overfitting: making our model too complicated, leading to higher training set accuracy but low testing set accuracy (an error due to high variance). Does not generalize well to new data.

K-Fold Cross Validation: splitting our data k time into training and testing sets to train and evaluate our model on different data combinations, so we don't lose some of the data.

Learning Curves: In a model that is underfitting, the training and CV error converge at a high error. In a model that is overfitting, the training and CV error do not converge. In a good model, the training and CV error will converge at a low error point.

Grid Search: combines all specified combinations of hyperparameters to find the best model. We then pick the one with the highest F1 score as our model and use a testing set to verify.

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.datasets import load_diabetes
4  from sklearn.model_selection import train_test_split, RandomizedSearchCV
5  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6  from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
7  import matplotlib.pyplot as plt
8  from sklearn.svm import SVC
9  import seaborn as sns
10 sns.set(style="ticks")
11
12 diabetes = pd.read_csv('diabetes.csv')
13
14 # Data Analysis
15 diabetes.describe() # see summary statistics for each column
16 diabetes.isna().any() # (0) - find number of NA values (by column)
17 diabetes['Outcome'].sum()/len(diabetes) # (0.35) - proportion of diabetes outcomes
18 diabetes['Age'].value_counts().hist() # right-skewed data
19 diabetes['Glucose'].hist() # approx normal distribution
20 sns.heatmap(diabetes.corr(), annot=True, cmap="YlGnBu"); # correlation heatmap
21
22 y = diabetes['Outcome']
23 X = diabetes.drop(columns='Outcome', axis=1)
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
25                                                    random_state=42)
26
27 # Build a RandomForestClassifier and use RandomizedSearch to find best params
28 clf_rf = RandomForestClassifier()
29
30 param_dist = {"max_depth": [3, None],
31              "n_estimators": list(range(10, 200)),
32              "max_features": list(range(1, X_test.shape[1]+1)),
33              "min_samples_split": list(range(2, 11)),
34              "min_samples_leaf": list(range(1, 11)),
35              "bootstrap": [True, False],
36              "criterion": ["gini", "entropy"]}
37
38 random_search = RandomizedSearchCV(clf_rf, param_distributions=param_dist)
39 random_search.fit(X_train, y_train)
40
41 rf_preds = random_search.best_estimator_.predict(X_test)
42 print('F1 Score: {}'.format(f1_score(y_test, rf_preds))) # 0.6847
```

```

1 # Build a AdaBoostClassifier and use RandomizedSearch to find best params
2 ada = AdaBoostClassifier()
3
4 parameters = {'n_estimators': list(range(50, 275, 25)),
5               'learning_rate': [0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 1, 2, 5, 10, 15]}
6
7 ada_random = RandomizedSearchCV(ada, param_distributions=parameters)
8 ada_random.fit(X_train, y_train)
9
10 ada_preds = ada_random.best_estimator_.predict(X_test)
11 print('F1 Score: {}'.format(f1_score(y_test, ada_preds))) # 0.6729
12
13 # Plotting feature importance
14 features = diabetes.columns[:-1]
15 importances = random_search.best_estimator_.feature_importances_
16 indices = np.argsort(importances)
17
18 plt.title('Feature Importance (Random Forest Classifier)')
19 plt.barh(range(len(indices)), importances[indices], color='b', align='center')
20 plt.yticks(range(len(indices)), [features[i] for i in indices])
21 plt.xlabel('Relative Importance')
22 plt.show() # Glucose, BMI, and Age were the top 3 features

```

Note that these features ranking were very similar to the heatmap that was made (except for pregnancies, which turned out to be the least influential feature).

In this case study, we looked at predicting diabetes for 768 patients. There was a reasonable amount of class imbalance with just under 35% of patients having diabetes. There were no missing data, and initial looks at the data showed it would be difficult to separate patients with diabetes from those that did not have diabetes.

Two advanced modeling techniques were used to predict whether or not a patient has diabetes. The most successful of these techniques proved to be the RandomForestClassifier, which had the following metrics:

Accuracy score for random forest : 0.7727272727272727

Precision score random forest : 0.6785714285714286

Recall score random forest : 0.6909090909090909

F1 score random forest : 0.6846846846846847

Based on the initial look at the data, it is unsurprising that Glucose, BMI, and Age were important in understanding if a patient has diabetes. It was surprising seeing that pregnancy looked to be correlated in the initial heatmap, but turned out to be the least important feature in our final model. This could have been due to its high correlation with the Age feature.

2 Deep Learning

2.1 Neural Networks

NOTE: refer back to Perceptron Algorithm section (1.2) for basic building blocks.

Error Function: tells us the distance we are from the points, assigning more weight to the misclassified points and less to the correctly classified points. This should be continuous (not discrete).

Continuous Predictions: we now want to classify points based on probability since we are working in the continuous world. Using the *Sigmoid Function*, we classify all points with probability < 0.5 as 0 and all points with probability > 0.5 as 1.

Sigmoid Function (Binary): we now classify points as $\hat{y} = \sigma(Wx + b)$. The formula is $\frac{1}{1+e^{-x}}$ and will give us a probability, which we classify based on a boundary line of 0.5.

Softmax Function (Multi-class): We apply the softmax function, $\frac{e^{z_i}}{e^{z_1} + \dots + e^{z_n}}$, to all of the linear function scores ($Wx + b$) to find the probability for each class.

```
1 def softmax(L):
2     # Takes a list of numbers and returns the softmax values for each
3     p_class = []
4     for z in L:
5         p_class.append(np.exp(z)/np.sum(np.exp(L)))
6     return p_class
7
```

One-Hot Encoding: for labels, create a column for each label option and assign a 1 if it is the corresponding label, otherwise assign it a 0. In general, we will have n columns for n labels.

Cross-Entropy: we want to maximize our probabilities that each point is correctly classified, which will minimize the error function. We sum across all of our probabilities, taking the $-\log()$ of each. If our point probability is near 1, it will have a small value. If our probability is near 0, it will have a larger value. This means a good model will give us a low cross-entropy. Points that are correctly classified have small errors, and points that are misclassified will have large errors.

```
1 def cross_entropy(Y, P):
2     # Y are our labels, P is our probabilities
3     Y = np.float_(Y)
4     P = np.float_(P)
5     return -np.sum(Y*np.log(P) + (1-Y)*np.log(1-P))
6
```

Binary Cross-Entropy: $E(W, b) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$ where $\hat{y}_i = \sigma(Wx^{(i)} + b)$

Multi-Class Cross-Entropy: $E(W, b) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y_{ij} \ln(\hat{y}_{ij})$ where $\hat{y}_{ij} = \sigma(Wx^{(ij)} + b)$

Gradient Descent: We want to minimize of size of $E(W, b)$ by taking steps downwards, which we can do by taking $\nabla E = [\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n}, \frac{\partial E}{\partial b}] = -(y - \hat{y})(x_1, \dots, x_n, 1)$ and then taking a step in the direction $-\alpha \nabla E$ since we want to move down scaled by a learning rate α . We then updated parameters as follows: $w'_i = w_i - \alpha \frac{\partial E}{\partial w_i} \implies w'_i = w_i + \alpha(y - \hat{y})x_i$ and $b' = b - \alpha \frac{\partial E}{\partial b} \implies b' = b + \alpha(y - \hat{y})$. We repeat this for a certain number of epochs until we minimize error.

```
1 # Activation (sigmoid) function
2 def sigmoid(x):
3     return 1/(1+np.exp(-x))
4
5 # Output (prediction) formula
6 def output_formula(features, weights, bias):
7     return sigmoid(np.dot(features, weights) + bias)
```

```

1 # Error (log-loss) formula
2 def error_formula(y, output):
3     return -y*np.log(output) - (1-y)*np.log(1-output)
4
5 # Gradient descent step
6 def update_weights(x, y, weights, bias, learnrate):
7     output = y-output_formula(x, weights, bias)
8     weights = weights + learnrate*output*x
9     bias = bias + learnrate*output
10    return weights, bias
11

```

General Neural Network: We have our input layer (x_1, \dots, x_n), our hidden layer (set of linear models created from the input), and our output layer (a combination of linear models to create a non-linear model). If we have n inputs, then our output will live in n dimensions.

Deep Neural Networks: We have an input layer, our linear models combine to create non-linear models in the hidden layers, and we combine these to output an advance non-linear model.

Multi-Class Networks: If we have n classes, then we have the same general architecture, but our output will have n nodes that each output a probability for the class. We then take the highest probability as the class label.

Feedforward: the process neural networks use to turn the input into an output. We take a vector of input features, and apply a sequence of linear models and sigmoid functions, which then output a probability of which label to assign to the prediction. Say we have a network with 1 input layer, 2 hidden layers, and an output layer. We take the output of the previous layer and use it as in input for our sigmoid of the current layer. We calculate \hat{y} as: $\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$

Error Function: $E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$ where \hat{y}_i is defined above.

Backpropagation: similar to gradient descent, we use this to update weights and get better predictions. However, keras handles this for us for the most part. To know: error term = $(y - \hat{y})\sigma'(x)$

2.2 Implementing Gradient Descent (NumPy)

Mean Square Error: the mean of the squares of the differences between the predictions and the labels.

Sum of Squared Errors (SSE): $E = \frac{1}{2} \sum_{\mu} \sum_j [y_j^{\mu} - \hat{y}_j^{\mu}]^2$ where variable j represents the output units of the network. The other sum over μ is a sum over all the data points. Adding these two errors and dividing by 2 gives us SSE.

Momentum: used to avoid ending up in a local minima instead of a global minima when performing gradient descent.

Weight Matrix: Each row in the matrix will correspond to the weights leading out of a single input unit (i), and each column will correspond to the weights leading in to a single hidden unit (j). So for a network with 3 inputs and 2 hidden layers, we will have a 3x2 matrix of weights, which we index at w_{ij} .

Input matrix: we need to be sure our dimensions match. The columns in the input must match the rows in the hidden layer.

Hidden Layer Input: we take the dot product of the input vector with the corresponding column of the weight matrix for the correct hidden layer unit.

```

1  def sigmoid(x):
2      return 1/(1+np.exp(-x))
3
4  # Network size
5  N_input = 4
6  N_hidden = 3
7  N_output = 2
8
9  np.random.seed(42)
10 X = np.random.randn(4)
11
12 weights_input_to_hidden = np.random.normal(0, scale=0.1, size=(N_input, N_hidden))
13 weights_hidden_to_output = np.random.normal(0, scale=0.1, size=(N_hidden, N_output))
14
15 hidden_layer_in = np.dot(X, weights_input_to_hidden) # 1x4*4x3 = 1x3
16 hidden_layer_out = sigmoid(hidden_layer_in) # 1x3 matrix passed through sigmoid
17
18 output_layer_in = np.dot(hidden_layer_out, weights_hidden_to_output) # 1x3*3x2 = 1x2
19 output_layer_out = sigmoid(output_layer_in) # 1x2 matrix passed through sigmoid
20

```

2.3 Training Neural Networks

Underfitting: an error due to high bias (too simple), low accuracy on both training and test sets.

Overfitting: an error due to high variance (too specific), low test accuracy but high training accuracy.

Model Choice: tend to the side of more complicated models, and apply techniques to prevent overfitting.

Model Complexity Graph: plots number of epochs against error, shows us when we start to overfit data as our testing error goes from decreasing to increasing when our training error is still decreasing.

Early Stopping: perform gradient descent until the testing error stops decreasing and starts increasing, at which point we stop training.

Regularization: large coefficients leads to our model overfitting, so we want to penalize large weights. We can redefine our error function using L1 or L2 regularization.

L1 Regularization: $E = -\frac{1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(|w_1| + \dots + |w_n|)$

When to use L1: using L1 we end up with sparse vectors (0,1,1,0,0,...) and small weights go to 0. If we want to reduce the number of weights and choose the best features for our model, this is a good option.

L2 Regularization: $E = -\frac{1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(w_2^2 + \dots + w_n^2)$

When to use L2: using L2 keeps all weights even if they are small. Normally this gives better results when training a model (smaller error), so this is the most common.

Dropout: some nodes have larger weights and more influence in our models training, so we choose a number of nodes to turn off randomly when training for each epoch (different each time). We assign a probability to each node of it being dropped when training.

Random Restart: when performing gradient descent, we can get stuck in a local minima instead of a global one. To avoid this, we can start from a few different random place and perform gradient descent from all of them (higher probability of finding global minima).

ReLU Activation: rectified linear unit, returns 0 if x is negative, otherwise returns x. The derivative is 1 if the number is positive, otherwise it is 0.

Batch Gradient Descent: we run all of our data through the network for each epoch, calculate the gradient for all points, and move in the given direction.

Stochastic Gradient Descent: take small subsets (batches) of our data and run them through the network, calculate the gradient for the subset of points and move in the given direction. Repeat this process for each batch of our data until we have gone through them all (this is considered one epoch).

Learning Rate: a large learning rate means big steps, which could go over our minima. A small learning rate means small steps, but longer training time. We can use a *decreasing learning rate* to change over time as we get closer to a minima.

Momentum: we want to take the average of the last few steps to determine the next step size. We multiply this by a constant, β , between 0 and 1 so more recent steps have more influence than farther back steps. $STEP(n) = STEP(n) + \beta STEP(n-1) + \beta^2 STEP(n-2) + \dots$

2.4 Deep Learning with PyTorch

2.4.1 Building Neural Networks

Tensors: the fundamental data structure for neural networks, similar to vectors and matrices.

Reshaping Tensors: we need input columns to match weights row. We can use `.view(r,c)` to do this.

NumPy to Torch: we often prepare data with NumPy and then want to convert it into tensors. The memory is shared between the arrays, so changing one will result in the other changing as well.

```
1  a = np.random.rand(4,3)
2  b = torch.from_numpy(a)
3  b.mul_(2)
4  # Note that now 'a' has been multiplied by 2 as well
5
```

MNIST Dataset: a collection of greyscale handwritten digits from 0-9. We read them in with a batch size of 64, and the trainloader contains batches of 64 images that we will get during each iteration (they will also be shuffled).

```
1  import numpy as np
2  import torch
3  import matplotlib.pyplot as plt
4  %matplotlib inline
5  %config InlineBackend.figure_format = 'retina'
6  from torchvision import datasets, transforms
7
8  # Define a transform to normalize the data
9  transform = transforms.Compose([transforms.ToTensor(),
10                                transforms.Normalize((0.5,), (0.5,))])
11
12 # Download and load the training data
13 trainset = datasets.MNIST('~/.pytorch/MNIST_data/', download=True, train=True,
14                           transform=transform)
15 trainloader = torch.utils.data.DataLoader(trainset, batch_size=64, shuffle=True)
16
17 dataiter = iter(trainloader) # grab first batch, used only to inspect some images
18 images, labels = dataiter.next()
19 print(images.shape) # torch.Size([64, 1, 28, 28]) [batch, color, height, width]
20
```


PyTorch Functional Module: we can define a model in PyTorch by using the functional library. We will build a 3 layer NN that uses a ReLU activation function and outputs to a softmax. We also want to initialize our weights and biases (this is done automatically, but we can also do it manually with custom values).

```
1  import torch.nn.functional as F
2
3  class Network(nn.Module): # subclass of nn.Module
4      def __init__(self):
5          super().__init__() # defines architecture
6          # Define input layers (Wx+b) for two hidden and one output layer
7          self.hidden1 = nn.Linear(784, 128) # initial input
8          self.hidden2 = nn.Linear(128, 64) # input from first hidden layer
9          self.output = nn.Linear(64, 10) # Output layer, 10 units - one for each digit
10
11      def forward(self, x):
12          x = F.relu(self.hidden1(x)) # Hidden layer1 with relu activation
13          x = F.relu(self.hidden2(x)) # Hidden layer2 with relu activation
14          x = F.softmax(self.output(x), dim=1) # Output layer with softmax activation
15          return x
16
17  model = Network()
18  model
19  # Network(
20  # (hidden1): Linear(in_features=784, out_features=128, bias=True)
21  # (hidden2): Linear(in_features=128, out_features=64, bias=True)
22  # (output): Linear(in_features=64, out_features=10, bias=True)
23  # )
24
25  # Set biases to all zeros (note that this is done inplace)
26  model.hidden1.bias.data.fill_(0)
27  # sample from random normal with standard dev = 0.01 (note that this is done inplace)
28  model.hidden1.weight.data.normal_(std=0.01)
29
```

PyTorch Sequential Module: lets build the same model above, using the sequential module.

```
1  from torch import nn
2
3  # Hyperparameters for our network
4  input_size = 784
5  hidden_sizes = [128, 64]
6  output_size = 10
7
8  # Build a feed-forward network
9  model = nn.Sequential(nn.Linear(input_size, hidden_sizes[0]),
10                        nn.ReLU(),
11                        nn.Linear(hidden_sizes[0], hidden_sizes[1]),
12                        nn.ReLU(),
13                        nn.Linear(hidden_sizes[1], output_size),
14                        nn.Softmax(dim=1))
15
16  # Forward pass through the network to classify single image
17  images, labels = next(iter(trainloader))
18  images.resize_(images.shape[0], 1, 784)
19  ps = model.forward(images[0,:])
20
```

2.4.2 Training Neural Networks

Loss function (cost): a measure of our prediction error after a forward pass. We then pass this through back propagation to update our weights using gradient descent. We want to use *cross entropy* as our loss for multi-class classification, but we will do so by using LogSoftmax and NLLLoss (negative log likelihood loss) together.

Autograd: a module used to automatically calculate gradient of tensors for backpropagation with respect to our loss function. By default when we create a NN, all parameters are initialized to require gradients (but you can manually turn them off if needed).

Optim: used to update the weights with the gradients, needs model parameters and learning rates. Note that we need to zero the gradients on each training pass otherwise you retain gradient from previous steps and weights do not update properly.

General Process (loop):

1. Make a forward pass through the network.
2. Use the network output to calculate the loss.
3. Perform a backward pass through the network to calculate gradients.
4. Take a step with the optimizer to update the weights.

```
1  from torch import optim
2
3  # Build Sequential model
4  model = nn.Sequential(nn.Linear(784, 128),
5                        nn.ReLU(),
6                        nn.Linear(128, 64),
7                        nn.ReLU(),
8                        nn.Linear(64, 10),
9                        nn.LogSoftmax(dim=1)) # sum across the columns
10
11 criterion = nn.NLLLoss() # Neg Log Likelihood Loss
12 optimizer = optim.SGD(model.parameters(), lr=0.003) # Stochastic Gradient Descent
13
14 epochs = 5
15 for e in range(epochs): # train for 5 steps
16     running_loss = 0
17     for images, labels in trainloader:
18         images = images.view(images.shape[0], -1) # Flatten images to a 784 long vector
19
20         optimizer.zero_grad() # clear any previous gradients
21         output = model.forward(images) # run images through model to get yhat
22         loss = criterion(output, labels) # calculate loss
23         loss.backward() # perform back pass & calculate gradients
24         optimizer.step() # move in the direction of gradient
25
26         running_loss += loss.item()
27     else:
28         print(f"Training loss: {running_loss/len(trainloader)}") # 1.89 -> 0.39
29
30 # Check model predictions
31 images, labels = next(iter(trainloader))
32 img = images[0].view(1, 784)
33 with torch.no_grad(): # Turn off gradients to speed up this part
34     logits = model.forward(img) # get
35
36 # Output of the network are logits, need to take softmax for probabilities
37 ps = F.softmax(logits, dim=1) # class probabilities
38 print(ps) # tensor of size [1,10] containing probability for each class
39
```

2.4.3 Validation and Dropout

topk: this returns the k highest values ($k=1$ means most likely class). This method returns a tuple of the top- k probability and the top- k indices. We can then compare these to their corresponding labels to see how many were correctly predicted.

Accuracy: we can take the mean of the comparison performed above between the predictions and labels to find out validation accuracy. We also turned the gradients off when doing this to speed up the process.

Overfitting: our training loss continues to decrease but our validation does not (could increase).

Dropout: regularization where we randomly drop input units so information is shared between the weights (increases ability to generalize to new data). However, when we are doing validation/testing we want to turn off dropout so we have all of our units.

```
1  import torch
2  from torchvision import datasets, transforms
3  from torch import nn, optim
4  import torch.nn.functional as F
5  import matplotlib.pyplot as plt
6
7  class Classifier(nn.Module): # Define a NN using dropout
8      def __init__(self):
9          super().__init__()
10         self.fc1 = nn.Linear(784, 256)
11         self.fc2 = nn.Linear(256, 128)
12         self.fc3 = nn.Linear(128, 64)
13         self.fc4 = nn.Linear(64, 10)
14         self.dropout = nn.Dropout(p=0.2)
15
16     def forward(self, x):
17         x = x.view(x.shape[0], -1) # flatten tensor
18         x = self.dropout(F.relu(self.fc1(x))) # hidden layer
19         x = self.dropout(F.relu(self.fc2(x))) # hidden layer
20         x = self.dropout(F.relu(self.fc3(x))) # hidden layer
21         x = F.log_softmax(self.fc4(x), dim=1) # output layer
22         return x
23
24 model = Classifier()
25 criterion = nn.NLLLoss()
26 optimizer = optim.Adam(model.parameters(), lr=0.003)
27
28 train_losses, test_losses = [], []
29 epochs = 10
30 for e in range(epochs):
31     running_loss = 0
32     for images, labels in trainloader: # training loop
33         optimizer.zero_grad() # clear all previous gradients
34         output = model.forward(images) # log probabilities
35         loss = criterion(output, labels) # calculate loss
36         loss.backward() # perform back prop and find gradient
37         optimizer.step() # take a step in direction of gradient
38         running_loss += loss.item() # add loss for each batch
39     else:
40         test_loss = 0
41         accuracy = 0
42         with torch.no_grad(): # turn off gradients
43             model.eval() # turn off dropout
44             for images, labels in testloader: # validation loop
45                 log_ps = model.forward(images) # log probabilities
46                 test_loss += criterion(log_ps, labels) # add loss for each batch
47
```

```

1     ps = torch.exp(log_ps) # get probabilities
2     top_p, top_class = ps.topk(k=1, dim=1) # get predicted class
3     equals = (top_class == labels.view(*top_class.shape)) # correct predictions
4     accuracy += torch.mean(equals.type(torch.FloatTensor)) # total correct
5
6     model.train() # turn back on dropout
7     train_losses.append(running_loss/len(trainloader)) # train loss for current epoch
8     test_losses.append(test_loss/len(testloader)) # test loss for current epoch
9
10    print("Epoch: {}/{}.. ".format(e+1, epochs),
11          "Training Loss: {:.3f}.. ".format(running_loss/len(trainloader)),
12          "Test Loss: {:.3f}.. ".format(test_loss/len(testloader)),
13          "Test Accuracy: {:.3f}".format(accuracy/len(testloader)))
14
15    # Inference (predict on new data)
16    model.eval() # turn off dropout
17    images, labels = next(iter(testloader)) # get image and label
18    img = images[0].view(1, 784) # Convert 2D image to 1D vector
19
20    with torch.no_grad(): # Calculate the class probabilities
21        output = model.forward(img) # log probabilities
22    ps = torch.exp(output) # tensor of probabilities
23

```

2.4.4 Saving and Loading Models

state_dict: after training a model, we want to save it to load later and either use to make predictions or continue training on. This way we don't have to retrain every time we load a workspace. We save a state_dict which contains the weights and bias matrices for each of our layers.

Saving and Loading: we need to rebuild the model exactly as it was when we load in a state_dict. This means saving in input size, output size, hidden layers, and the weights/bias.

```

1     checkpoint = {'input_size': 784,
2                  'output_size': 10,
3                  'hidden_layers': [each.out_features for each in model.hidden_layers],
4                  'state_dict': model.state_dict()}
5
6     torch.save(checkpoint, 'checkpoint.pth')
7
8     def load_checkpoint(filepath):
9         checkpoint = torch.load(filepath)
10        model = ... # set up correct model architecture
11        model.load_state_dict(checkpoint['state_dict'])
12        return model
13
14    model = load_checkpoint('checkpoint.pth')
15

```

2.4.5 Loading Image Data

ImageFolder: we want to have a directory that contains all of our images, separated by train/test folders within, and separated by classes within the train/test directories.

Transform: we can augment images when reading them in by using the transform method. Note that we only want to do this to the training data and now the testing or validation data sets. The training data will just get the normal transform methods.

Data Loaders: takes a dataset and returns batches of images/labels (can also shuffle data after each epoch). You have to convert it to an iterator and call `next()` since it is a generator.

```
1 data_dir = 'Cat_Dog_data' # directory containing train/test folders
2
3 train_transforms = transforms.Compose([transforms.RandomRotation(30), # augment
4                                       transforms.RandomHorizontalFlip(), # augment
5                                       transforms.Resize(255), # 255x255 image
6                                       transforms.CenterCrop(224),
7                                       transforms.ToTensor() # convert to tensor
8                                       ])
9
10 test_transforms = transforms.Compose([transforms.Resize(255), # 255x255 image
11                                      transforms.CenterCrop(224),
12                                      transforms.ToTensor() # convert to tensor
13                                      ])
14 # Note that ToTensor automatically converts color range to [0,1]
15
16 # Load datasets from folders, pass transformations as they are loaded
17 train_data = datasets.ImageFolder(data_dir + '/train', transform=train_transforms)
18 test_data = datasets.ImageFolder(data_dir + '/test', transform=test_transforms)
19
20 # Load into generator with batch size 32 (shuffle training data)
21 trainloader = torch.utils.data.DataLoader(train_data, batch_size=32, shuffle=True)
22 testloader = torch.utils.data.DataLoader(test_data, batch_size=32)
23
```

2.4.6 Transfer Learning

Case 1 (small data set, similar data): we want to drop the output layer and add our own fully connected layer to match the number of classes. We freeze all weights of the other layers but initialize random weights for new output layer, training the network only on these new output layer weights. This prevents overfitting since the data is small.

Case 2 (small data set, different data): we want to only keep the early layers (not the high level features since they are different data) and add a new fully connected layer to match the number of classes. We freeze the early layers weights and train the network only for the new output layer. This prevents overfitting since the data is small.

Case 3 (large data set, similar data): drop the output layer and add a fully connected layer with the number of classes for our data. Since we have a large dataset, we retrain all of the weights in the network since overfitting is less of a concern (fine-tuning the weights).

Case 4 (large data set, different data): we use the same approach as case 3 above. However, if this does not give us a successful model we can randomly initialize the weights and train the network from scratch.

Matching Data: most pretrained models use an image size of 224x224 as input. They also require you to normalize the color channels in the same way that they were trained, with mean [0.485, 0.456, 0.406] and std [0.229, 0.224, 0.225].

```
1 import torch
2 from torch import nn
3 from torch import optim
4 import torch.nn.functional as F
5 from torchvision import datasets, transforms, models
6 from collections import OrderedDict
7
8 data_dir = 'Cat_Dog_data'
```

```

1  train_transforms = transforms.Compose([transforms.RandomRotation(30),
2                                         transforms.RandomResizedCrop(224),
3                                         transforms.RandomHorizontalFlip(),
4                                         transforms.ToTensor(),
5                                         transforms.Normalize([0.485, 0.456, 0.406],
6                                                                [0.229, 0.224, 0.225])])
7
8  test_transforms = transforms.Compose([transforms.Resize(255),
9                                         transforms.CenterCrop(224),
10                                        transforms.ToTensor(),
11                                        transforms.Normalize([0.485, 0.456, 0.406],
12                                                            [0.229, 0.224, 0.225])])
13
14  train_data = datasets.ImageFolder(data_dir + '/train', transform=train_transforms)
15  test_data = datasets.ImageFolder(data_dir + '/test', transform=test_transforms)
16
17  trainloader = torch.utils.data.DataLoader(train_data, batch_size=64, shuffle=True)
18  testloader = torch.utils.data.DataLoader(test_data, batch_size=64)
19
20  model = models.densenet121(pretrained=True) # import model with pretrained weights
21
22  # Use GPU if it's available (otherwise use a CPU)
23  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
24
25  for param in model.parameters(): # freeze all layers
26      param.requires_grad = False # turn off gradients
27
28  model.classifier = nn.Sequential(OrderedDict([ # replace old models classifier
29                                             ('fc1', nn.Linear(1024, 500)),
30                                             ('relu', nn.ReLU()),
31                                             ('fc2', nn.Linear(500, 2)),
32                                             ('output', nn.LogSoftmax(dim=1)) ]))
33
34  criterion = nn.NLLLoss() # define loss function
35  optimizer = optim.Adam(model.classifier.parameters(), lr=0.003)
36
37  model.to(device); # move model to either GPU or CPU
38
39  epochs = 1 # train for one epoch
40  steps = 0
41  running_loss = 0 # training loss
42  print_every = 5
43  for epoch in range(epochs):
44      for inputs, labels in trainloader: # training loop
45          steps += 1
46          # Move input and label tensors to the default device
47          inputs, labels = inputs.to(device), labels.to(device)
48
49          optimizer.zero_grad() # clear previous gradients
50          logps = model.forward(inputs) # log probabilities
51          loss = criterion(logps, labels) # calculate loss
52          loss.backward() # back prop gradient descent
53          optimizer.step()
54
55          running_loss += loss.item()
56
57      if steps % print_every == 0: # every 5 steps (validation loop)
58          test_loss = 0
59          accuracy = 0
60          model.eval() # turn off dropout

```

```

1  with torch.no_grad(): # turn off gradient descent
2
3      for inputs, labels in testloader: # loop through validation
4          inputs, labels = inputs.to(device), labels.to(device)
5          logps = model.forward(inputs) # log probabilities
6          batch_loss = criterion(logps, labels)
7
8          test_loss += batch_loss.item()
9
10         ps = torch.exp(logps) # probabilities
11         top_p, top_class = ps.topk(1, dim=1) # get top class
12         equals = top_class == labels.view(*top_class.shape)
13         accuracy += torch.mean(equals.type(torch.FloatTensor)).item()
14
15     print(f"Epoch {epoch+1}/{epochs}.. "
16           f"Train loss: {running_loss/print_every:.3f}.. "
17           f"Test loss: {test_loss/len(testloader):.3f}.. "
18           f"Test accuracy: {accuracy/len(testloader):.3f}")
19     running_loss = 0 # reset training loss
20     model.train() # turn dropout back on
21

```

3 Unsupervised Learning

3.1 Clustering (KMeans)

Prior Knowledge: when choosing k , we may have an idea of how many clusters we want to create based off of a question that have been presented.

Decision Method: when choosing k , we don't know how many clusters we want our data to be in.

Elbow Method: create a plot of the number of clusters (on the x-axis) vs. the average distance of the center of the cluster to each point (on the y-axis). We want to chose the k where we stop decreasing our distance by a large amount (no substantial change from increasing k).

```
1 scores = []
2 centers = list(range(1,11))
3 for k in centers: # fit k from 1 to 10
4     model = KMeans(k).fit(data)
5     scores.append(abs(model.score(data))) # distance from center
6
7 plt.plot(centers, scores, linestyle='--', marker='o')
8 plt.ylabel('Score (Avg Distance)')
9 plt.xlabel('K-value (Centers)')
10 plt.title('Scree Plot'); # we see 4 is the correct k
11
```

Algorithm: randomly place k centroids in your data. Then assign each point to the closest centroid and move it to the center of the points assigned (repeat this until the centroid stops moving).

Optimal Clusters: to make sure you have the best clusters, KMeans will be performed multiple times to chose the best set of clusters that have the smallest average distance from point to centroid.

Scalers: *StandardScaler* will transform data to have a mean of 0 and standard deviation of 1.

MinMaxScaler will transform the data to be in the range [0,1] with respect to the max value. Both will result in different clusters, but *StandardScaler* is the most common for clustering algorithms.

3.2 Hierarchical and Density Based Clustering

3.2.1 Hierarchical Clustering

Single-Link Clustering: we specify 1 to n (number of examples) clusters that we want to make. It looks at the closest points to one another in order to make clusters. A warning is that this method can eat up a lot of points and cause huge imbalances.

Dendrograms: shows us how the hierarchical cluster is separated. No matter how many dimensions are data is in, we can also produce this graph to visualize the data.

Complete-Link: similar to single link, but the distance measure is now the farthest point from a cluster instead (this produces compact clusters).

Average-Link: similar to single link, but the distance measure is between every point in the clusters and the average is taken from this value.

Ward's Method: (default method in sklearn). This involves finding a central point between two clusters and taking the squared distance between all points in the clusters and adding them together. You then subtract the variance of the two clusters to find the total distance. (minimizes variance).

Adjusted Rand Score: a score between -1 and 1, where 1 means two clusterings are identical of how they grouped the samples in a dataset.


```

1  from sklearn import datasets
2  from sklearn.cluster import AgglomerativeClustering
3  from sklearn.metrics import adjusted_rand_score
4  from sklearn import preprocessing
5  from scipy.cluster.hierarchy import linkage, dendrogram
6
7  iris = datasets.load_iris()
8  normalized_X = preprocessing.normalize(iris.data) # range [0,1]
9
10 ward = AgglomerativeClustering(n_clusters=3)
11 ward_pred = ward.fit_predict(normalized_X)
12
13 complete = AgglomerativeClustering(n_clusters=3, linkage="complete")
14 complete_pred = complete.fit_predict(normalized_X)
15
16 avg = AgglomerativeClustering(n_clusters=3, linkage="average")
17 avg_pred = avg.fit_predict(normalized_X)
18
19 ward_ar_score = adjusted_rand_score(iris.target, ward_pred) # 0.8857
20 complete_ar_score = adjusted_rand_score(iris.target, complete_pred) # 0.6444
21 avg_ar_score = adjusted_rand_score(iris.target, avg_pred) # 0.5584
22
23 linkage_type = 'ward'
24 linkage_matrix = linkage(normalized_X, linkage_type) # ward clustering
25 dendrogram(linkage_matrix)
26 plt.show();
27
28 # We can also use seaborn which performs clustering and graphs in one method
29 sns.clustermap(normalized_X, figsize=(12,18), method=linkage_type, cmap='viridis')
30 plt.show();
31

```

Advantages: hierarchical representation can be informative and provides ability to visualize, especially good when dataset contains real hierarchical relationships (evolutionary biology).

Disadvantages: sensitive to noise/outliers, computationally intensive with $O(N^2)$.

3.2.2 Density-Based Clustering

DBSCAN: groups together points that are closely packed together. It does this by a factor of ϵ (which determines search distance around a point) and *MinPts* (minimum number of points required to form a cluster). It labels points that don't fall into these clusters as noise. Points that have clusters around them are labeled *core points* and the points are them that make up the cluster are called *border points*.

```

1  from sklearn import cluster
2
3  epsilon = .75
4  min_samples = 10
5
6  dbscan = cluster.DBSCAN(eps=epsilon, min_samples=min_samples)
7  clustering_labels_2 = dbscan.fit_predict(dataset_2)
8

```

Advantages: don't need to specify number of clusters, flexibility in shape/size of clusters, able to deal with noise and outliers.

Disadvantages: border points that are reachable from two clusters, difficulty finding clusters of varying densities.

3.3 Gaussian Mixture Models and Cluster Validation

GMM Clustering: assumes that each cluster follows the normal distribution.

Multivariate Gaussian Distribution: more than one variable follows the normal distribution. We can plot a point in the middle of a cluster (the mean of the normal distributions) and each ellipse represents 1 standard deviation from the mean.

Expectation Maximization Algorithm: initialize k Gaussian distributions, soft-cluster the data (“expectation”), re-estimate the Gaussian’s (“maximization”), and evaluate log-likelihood to check for convergence. If no convergence then repeat from step 2.

Step 1: set mean and variance for each Gaussian distribution by running k -means and using clusters produced to initialize our parameters.

Step 2: find the membership likelihood (probability) that each points belongs to each clusters.

Step 3: find the new weighted values for mean and variance and updated the Gaussian distributions.

Step 4: estimate the log-likelihood, we want to maximize this value (converge to a high value).

Covariance Type: determine the shape of the clusters (sphere, full, etc.).

```
1  from sklearn import datasets, mixture
2
3  X = datasets.load_iris().data[:10]
4
5  gmm = mixture.GaussianMixture(n_components=3)
6  gmm.fit(X)
7  clustering = gmm.predict(X)
8
```

Advantages: soft-clustering (sample membership of multiple clusters) and flexibility of cluster shape.

Disadvantages: sensitive to initialization values, possible converge to a local optimum, and a slow convergence rate.

3.3.1 Cluster Analysis Process

Feature Selection: we may want to cluster based on a select few features (if we have 1000’s of dimensions we may only want to use a few features).

Feature Extraction: transforming the data to generate useful features (PCA).

Algorithm Selection: multiple models must be tried to find the best fit (no one model is the best). The *proximity measure* is used to calculate the distance between points (different for each problem type).

Cluster Validation: *indices*, evaluating how well a clustering turned out to be (score of the quality).

- **External Indices:** used if our data was originally labeled (Adjusted Rand Score).
- **Internal Indices:** measure the fit between data and clusters only using the data.
- **Relative Indices:** which of two clustering structures is the best.
- **Compactness:** measure of how close the elements of a cluster are to each other.
- **Separability:** measure of how far clusters are from each other.

Results Interpretation: determine what insights can be gotten from the clusters (assigning labels).

Silhouette Coefficient: a measure for *internal indices* that uses the average distance between a clusters points compared to the next closest clusters points. This can help us determine our value of k (the closer to 1 the better). We can also use this to compare the performance of different algorithms on a dataset. Works best for dense, circular clusters. Don’t use for DBSCAN.

make_blob: an sklearn method that creates a dataset of Gaussian blobs.

```
1  import matplotlib.pyplot as plt
2  from sklearn import cluster, datasets, mixture
3  from sklearn.cluster import KMeans
4  from sklearn.metrics import adjusted_rand_score
5  from sklearn.mixture import GaussianMixture
6  %matplotlib inline
7  import seaborn as sns
8
9  iris = sns.load_dataset("iris")
10 features = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width']
11
12 # PairGrid - plots every pair of features against each other in a subplot
13 g = sns.PairGrid(iris, hue="species", palette=sns.color_palette("cubehelix", 3),
14                 vars=features)
15 g.map(plt.scatter)
16 plt.show()
17
18 # Predict clusters using K-Means and add to dataframe
19 kmeans_iris = KMeans(n_clusters=3)
20 pred_kmeans_iris = kmeans_iris.fit_predict(iris[features])
21 iris['kmeans_pred'] = pred_kmeans_iris
22
23 # Predict using GMM and add to dataframe
24 gmm_iris = mixture.GaussianMixture(n_components=3).fit(iris[features])
25 pred_gmm_iris = gmm_iris.predict(iris[features])
26 iris['gmm_pred'] = pred_gmm_iris
27
28 # Find Adj Rand Score (since we have too many dimensions to visualize)
29 iris_kmeans_score = adjusted_rand_score(iris.species, iris.kmeans_pred) # 0.73
30 iris_gmm_score = adjusted_rand_score(iris.species, iris.gmm_pred) # 0.904 (best)
31
```

3.4 Dimensionality Reduction and PCA

PCA: taking the dataset and reducing it to only the parts that hold the most information.

Latent Features: features that are not directly in our dataset, but encompass features that already exist. Putting features of our data into broad categories to reduce dimensions.

Feature Selection (doc): finding a subset of the original data that is most relevant and useful.

- **Filter methods:** finding correlation between data (unsupervised) or labels (supervised) and filtering out the less useful ones. Methods used include Pearson's Correlation, LDA, and ANOVA.
- **Wrapper methods:** selecting features based on their impact on the model performance. Multiple models and cross validation are used, and the features that produce the best model are kept. Methods used include Forward Search, Backward Search, and Recursive Feature Elimination.

Feature Extraction: extracting new 'latent features' from original dataset. Keeps the most information possible from multiple features while also reducing dimensionality at a maximum. This is the goal of PCA, but can also be accomplished through ICA and Random Projection.

Principal Components: are linear combinations of the original features in a dataset that aim to retain the most information in the original data (similar to latent features). We want to find a projection of the data onto a smaller number of features, while preserving the maximum amount of information.

PCA Properties: (1) each component captures the largest amount of variance in the dataset (reduce the distance each point is from the component). (2) components are orthogonal to one another, meaning each when we have multiple components they must be 90deg from each other (independent features guarantees this).

```
1  from sklearn.decomposition import PCA
2  from sklearn.preprocessing import StandardScaler
3  from sklearn.ensemble import RandomForestClassifier
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import accuracy_score
6
7  def fit_rfc(X,y):
8      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33,
9          random_state=42)
10     clf = RandomForestClassifier(n_estimators=100, max_depth=None)
11     clf.fit(X_train, y_train)
12     y_preds = clf.predict(X_test)
13     acc = accuracy_score(y_test, y_preds)
14     return acc
15
16 def do_pca(n_components, data):
17     X = StandardScaler().fit_transform(data)
18     pca = PCA(n_components) # pca object
19     X_pca = pca.fit_transform(X) # new matrix of data
20     return pca, X_pca
21
22 def pca_optimize(X, y, n_components=100):
23     accs = []
24     for component in range(3, n_components):
25         pca, X_pca = do_pca(component, X)
26         accs.append(fit_rfc(X_pca, y))
27     return accs
28
29 def plot_pca_acc(acc):
30     plt.plot(range(3, len(acc)+3), acc)
31     plt.title('PCA N-Components Accuracy')
32     plt.ylabel('Accuracy')
33     plt.xlabel('N-Components')
34     return None
35
36 acc = pca_optimize(X,y,50)
37 plot_pca_acc(acc) # around 15 components we begin to level off around 0.91
```

Interpreting PCA: the PCA object we create has many methods we can use to interpret the new components.

Explained Variance: the *explained_variance_ratio_* method will tell us how much variance each components represents. Note that the first component always explains the most amount of variability. Summing up all of the variance among components tells use how much of the total variability the components capture. There are called *eigenvalues*

Component Weights: the *components_* method returns a vector that explains the weight given to each component. For an image, this is the weight given to each pixel (which tells use the most important aspect of the image). There are called *eigenvectors*.

```

1 df = pd.read_csv('./data/cars.csv')
2
3 # We want to components to explain 85% of variability of original datasets
4 for ncomp in range(3,30):
5     pca, X_pca = do_pca(ncomp, df)
6     if sum(pca.explained_variance_ratio_) > 0.85:
7         break
8
9 pca.n_components # 6
10 sum(pca.explained_variance_ratio_) # 0.8683
11
12 # Create a df and visualize results (see pca_helper_functions in Files folder)
13 pca_results(df, pca)
14

```

3.5 Random Projections and ICA

Random Projection: we multiply a dataset by a random matrix in order to reduce the number of features while still preserving the distance between the points. We use ϵ to determine the level of error that we allow the reduction to have. This is also more computationally efficient than PCA, so it works better with higher dimensional data (but is a decrease in projection quality).

```

1 from sklearn import random_projection
2
3 rp = random_projection.SparseRandomProjection()
4 new_X = rp.fit_transform(X) # new transformed dataset
5

```

Independent Component Analysis: assumes features are a mixture of independent sources and tries to isolate them. Note that we need as many observations as the original signals we are trying to separate.

FastICA Algorithm: assumes that our components are independent and do not follow a Gaussian distribution. We have X (the original recordings), which the product of multiplying A (mixing matrix) by S (the source signal). We can find the source single by $S = WX$ where $W = A^{-1}$ (the unmixing matrix). We want to find the best W to produce the original signals.

Algorithm Steps: (1) center, whiten X . (2) choose initial random weights matrix W . (3) estimate W , containing vectors. (4) decorrelate W to prevent any W_i from converging to same value. (5) repeat step 3 and 4 until converged.

```

1 import numpy as np
2 import wave
3
4 # Read in wave files
5 mix_1_wave = wave.open('ICA mix 1.wav', 'r')
6 mix_2_wave = wave.open('ICA mix 2.wav', 'r')
7 mix_3_wave = wave.open('ICA mix 3.wav', 'r')
8 mix_1_wave.getparams() # see details (channels, length, etc).
9
10 # Extract Raw Audio from Wav Files
11 signal_1_raw = mix_1_wave.readframes(-1)
12 signal_1 = np.fromstring(signal_1_raw, 'Int16')
13
14 signal_raw_2 = mix_2_wave.readframes(-1)
15 signal_2 = np.fromstring(signal_raw_2, 'Int16')
16
17 signal_raw_3 = mix_3_wave.readframes(-1)
18 signal_3 = np.fromstring(signal_raw_3, 'Int16')

```

```

1  import matplotlib.pyplot as plt
2  from sklearn.decomposition import FastICA
3  from scipy.io import wavfile
4
5  fs = mix_1_wave.getframerate()
6  timing = np.linspace(0, len(signal_1)/fs, num=len(signal_1))
7
8  # Plot first recording
9  plt.figure(figsize=(12,2))
10 plt.title('Recording 1')
11 plt.plot(timing, signal_1, c="#3ABFE7")
12 plt.ylim(-35000, 35000)
13 plt.show();
14
15 X = list(zip(signal_1, signal_2, signal_3)) #
16
17 ica = FastICA(n_components=3)
18 ica_result = ica.fit_transform(X)
19 ica_result.shape # (264515, 3)
20
21 # Seperate each signal into an array
22 result_signal_1 = ica_result[:,0]
23 result_signal_2 = ica_result[:,1]
24 result_signal_3 = ica_result[:,2]
25
26 # Convert to int, map the appropriate range (int16), and increase the volume
27 result_signal_1_int = np.int16(result_signal_1*32767*100)
28 result_signal_2_int = np.int16(result_signal_2*32767*100)
29 result_signal_3_int = np.int16(result_signal_3*32767*100)
30
31 # Write wave files
32 wavfile.write("result_signal_1.wav", fs, result_signal_1_int)
33 wavfile.write("result_signal_2.wav", fs, result_signal_2_int)
34 wavfile.write("result_signal_3.wav", fs, result_signal_3_int)
35

```