

W1 - Intro

Parameter-based, ML System

W2 - Linear regression

$$f(x) = \log(x_1, w_1) + w_2 x_2 + w_3$$

parameters linear, inputs non-linear

MSE (Mean squared error)

$$L(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$L(w, w_i) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_i x_i))^2$$

RSS (residual sum of squares)

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (\epsilon_i)^2$$

covariance

$$w^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{6xg}{6x^2}$$

$$\text{correlation coefficient } r_{xy} = \frac{6xg}{6x^2}$$

$$L(w^*, w_i) = 1 - \frac{6xg}{6x^2}$$

$$\Rightarrow \frac{MSE}{\sigma_y^2} = 1 - R^2$$

$$R^2 = 1 - \frac{MSE}{\sigma_y^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$R^2 \approx 1 \rightarrow$ model explains all variance of y

$R^2 \approx 0 \rightarrow$ doesn't explain any of ...

$$f_i = \sum_{j=0}^d w_j e^{jx_i} \quad \vec{f}(x) = [f_0(x), \dots, f_d(x)]$$

$$\hat{y} = \vec{w}^\top \vec{x} \quad \text{given } (x_1, y_1), \dots, (x_n, y_n)$$

$$\vec{I} = \begin{bmatrix} 1 & e^{-x_1} & e^{-2x_1} \\ 1 & e^{-x_2} & e^{-2x_2} \\ 1 & e^{-x_3} & e^{-2x_3} \\ 1 & e^{-x_4} & e^{-2x_4} \end{bmatrix}$$

$$J = C + Dx + E, E \sim N(0, \sigma^2)$$

training data from $y = A + Bx + E, R^2 \approx 1 (6=0)$

test data on $y = C + Dx + E, R^2 \approx 1 (6=0)$

W3 - Gradient descent, bias variance

$$L(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$$

$$w^{t+1} = w^t - \frac{\partial}{\partial w} \sum_{i=1}^n (-x_i)(y_i - \langle w^t, x_i \rangle)$$

$$i \quad x_0 \quad x_1 \quad y_1 \quad \langle -x_1, y_1 \rangle$$

$$0 \quad 1 \quad -7.49 \quad -0.919 \quad \langle 7.49, -15.81 \rangle$$

$$1 \quad 1 \quad 1.748 \quad -1.048 \quad \langle -10.69, -18.69 \rangle$$

$$2 \quad 1 \quad 3.212 \quad 1.9117 \quad \langle -38.17, -12.64 \rangle$$

$$w_0 = 3, w_1 = -5, w_2 = 0.1 \quad \langle -40.87, -57.04 \rangle$$

$$\hat{y}_i = w_0 x_0 + w_1 x_1 + w_2 x_2 = \langle 6.95, -11.74, -19.06 \rangle$$

$$W^{t+1} = [-3, -5] - \frac{0.1}{3} \left[\begin{array}{c} 40.87 \\ -57.04 \end{array} \right] = [-4.64, 0.13]$$

$$Momentum: update includes a vector v, that accumulates gradient of past steps. Each update is a linear combination of the gradient and the previous update. (Go faster if gradient keeps pointing in the same direction)$$

AdaGrad: per-parameter learning rate. Track per-parameter of gradient to normalize parameter update step. Weights with large gradient have smaller learning rate, weights with small gradient have larger learning rate. Take smaller steps in steep directions, take bigger steps where the gradient is flat.

RMSProp: Leaky AdaGrad. Uses EWMA to emphasize recent gradient magnitudes

Adam: Adaptive Moment Estimation. Uses ideas from momentum (first moment) and RMSProp (second moment) + plus bias correction.

Total 1600 samples, in one epoch:

how many samples to full batch / mini-batch / per epoch / per iteration / one sample

compute the gradient of loss for that sample

how many times to update weights (iterations)

$\vec{f}(x_t, w) \rightarrow f; t(x_t) \rightarrow t$

$$E[(y_i - \hat{y}_i)^2] = t - f^2 + \sigma_e^2$$

$$E[(t-f)^2 + \sigma_e^2] = (t-E[f])^2 + (E[f]^2) - E[(f)^2] + \sigma_e^2$$

$$(Bias)^2 \leq \text{variance over } D \leq \text{irreducible error}$$

$$\text{Var} = \frac{p}{n} \sigma_e^2 \quad \text{number of parameters} \quad \text{number of samples}$$

$$y = A + Bx + E, E \sim N(0, \sigma^2)$$

training data: $0 \leq R^2 \leq 1, R^2 = 1 (\sigma=0)$

test data: $R^2 \approx 1, R^2 \approx 1 (\sigma=0)$

W4 - Model selection and regularization

Hold-out validation

divide data into training, validation, test sets.

Learn on training set for each candidate model.

Measure error for all models on validation set.

Select model that minimizes error on validation set.

Evaluate that model on test set.

Problems: fitted model (and test error) varies a lot depending on sample selected for training and validation. Fewer samples available for estimating parameters. Especially bad for problems with small number of samples.

K-fold cross-validation

Divide data into K equal-sized pools (typically 5, 10)

For each, evaluate model. $K-1$ parts training, 1 part validation.

Average K validation scores and choose based on avg.

Leave-p-out CV

In each iteration, p validation points

Remaining n-p points are for training

Repeat for all possible sets of p validation points

K-fold CV - how to split?

No structure — shuffle split (arid accidental problems)

Group structure — keeps members of each group in either training set, or validation set, but not both.

Time series data — keeps validation data in the future, relative to training data

One standard error rule: use simplest model where mean error is within 1 SE of the minimum one.

Given data $X \cdot y$, compute score S_p , for model p on fold i

Compute avg (S_p), standard deviation σ_p , and standard error of score: $SEP = \frac{\sigma_p}{\sqrt{K-1}}$

For MSE, $p^* = \arg \min_p S_p, S_p = \sum_i (y_i - \hat{y}_i)^2$

For R2, $D^* = \arg \max_p S_p, S_p = \sum_i (y_i - \hat{y}_i)^2, S_p^* = \min_p \{S_p | \hat{y}_i \neq y_i\}$

Ridge regression ($L2$): $\hat{y}(w) = \sum_{j=1}^d w_j x_j$

LASSO regression ($L1$): $\hat{y}(w) = \sum_{j=1}^d w_j x_j$

Both penalize large w_j ; Both have parameter that controls levels of regularization; Intercept not in regularization sum, it depends on mean of y .

Ridge ($L2$)

minimizes $\|w\|_2^2$, minimal penalty for non-zero coeffs.

having penalties large coeffs

tends to make many "small" coeffs; Not for feature selection

LASSO ($L1$)

minimizes $\|w\|_1$, tends to make coeffs either 0 or larger (sparse)

does feature selection ($w_j \neq 0 \Rightarrow$ unselecting feature)

Ridge ($L2$)

$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d |w_j|^2 = \|A w - y\|^2 + \lambda \|w\|^2$

LASSO $J(w) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d |w_j| = \|A w - y\|^2 + \lambda \|w\|_1$

Greater λ , less complex model. \rightarrow bias, variance. Ridge \rightarrow coefficients smaller; LASSO \rightarrow more weights zero

COVID

```

df_tr_poly = df_tr.daysElapsed.values.reshape(-1,1)*np.arange(1, 5)
df_ts_poly = df_ts.daysElapsed.values.reshape(-1,1)*np.arange(1, 5)
reg_covid = LinearRegression().fit(df_tr_poly, df_tr.deathIncrease)
deathIncrease_fitted = reg_covid.predict(df_tr_poly)
deathIncrease_fitted = np.exp(df_tr_poly.deathIncrease, deathIncrease_fitted)
deathIncrease_fitted

```

```

def kfold_cv(deathIncrease, deathIncrease_fitted):
    nfold = 5
    dtest_list = np.arange(1,nmax+1)
    idm = len(dtest_list)
    idval = np.array([10, 18, 48, 50])
    nfold = len(idval)
    mse_val = np.zeros((nfold))
    r2_val = np.zeros((nfold))
    mse_tr = np.zeros((nfold))
    r2_tr = np.zeros((nfold))
    r2_val = np.zeros((nfold))
    r2_tr = np.zeros((nfold))

    for fold_idx in range(nfold):
        for dtest, idx in enumerate(dtest_list):
            x_train_kFold = df_tr_poly[:,idx]
            y_train_kFold = df_tr.deathIncrease.values[:,idx]
            x_val_kFold = df_tr_poly[idval[fold_idx]:idval[fold_idx]+1]
            y_val_kFold = df_tr.deathIncrease.values[idx:idx+10]
            for dtest, idx in enumerate(dtest_list):
                x_train_kTest = x_train_kFold[:,idx]
                y_train_kTest = y_train_kFold[:,idx]
                reg_dtest = LinearRegression().fit(x_train_kTest, y_train_kTest)
                y_hat = reg_dtest.predict(x_val_kFold)
                mse_val[fold_idx], isplit = metrics.mean_squared_error(y_val_kFold, y_hat)
                r2_val[fold_idx], isplit = metrics.r2_score(y_val_kFold, y_hat)
                r2_tr[fold_idx], isplit = metrics.r2_score(x_train_kTest, y_hat)
                mse_tr[fold_idx], isplit = metrics.mean_squared_error(y_train_kFold, y_hat_tr)
                r2_tr[fold_idx], isplit = metrics.r2_score(y_train_kFold, y_hat_tr)

```

Linear Regression K-fold

The outer loop can be used to divide the data into training and validation, but then we also need an inner loop to train and validate each model for this particular fold.

In this case, we want to evaluate polynomial models with different model orders from

$$d=1, \hat{y} = w_1x + w_0$$

$$d=10, \hat{y} = w_{10}x^{10} + w_9x^9 + \dots + w_0$$

We could do something like this:

```

nfold = 5
kf = KFold(n_splits=nfold, shuffle=True)
dmax = 10
dtest_list = np.arange(1,nmax+1)
idm = len(dtest_list)
idval = np.array([10, 18, 48, 50])
mse_val = np.zeros((nfold))
r2_val = np.zeros((nfold))
mse_tr = np.zeros((nfold))
r2_tr = np.zeros((nfold))

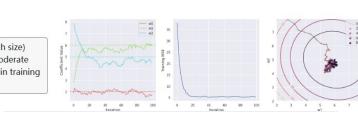
# Loop over the folds
# The first loop variable tells us how many out of nfold folds we have gone through
# The second loop variable tells us which fold we are splitting the data
# For this iteration, we are splitting kf.get_n_splits(kf.get_n_splits())
# These are the indices for the training and validation indices
# This is the iteration of the k folds
for idk, idx in enumerate(dtest_list):
    x_train_kFold = x_train[:,idx]
    y_train_kFold = y_train[:,idx]
    x_val_kFold = x_train[idval[0]:idval[-1],idx]
    y_val_kFold = y_train[idval[0]:idval[-1],idx]
    for dtest, idx in enumerate(dtest_list):
        x_train_kTest = x_train_kFold[:,idx]
        y_train_kTest = y_train_kFold[:,idx]
        reg_dtest = LinearRegression().fit(x_train_kTest, y_train_kTest)
        y_hat = reg_dtest.predict(x_val_kFold)
        mse_val[0], isplit = metrics.mean_squared_error(y_val_kFold, y_hat)
        r2_val[0], isplit = metrics.r2_score(y_val_kFold, y_hat)
        r2_tr[0], isplit = metrics.r2_score(x_train_kTest, y_hat)
        mse_tr[0], isplit = metrics.mean_squared_error(y_train_kFold, y_hat_tr)
        r2_tr[0], isplit = metrics.r2_score(y_train_kFold, y_hat_tr)

        # Now
        idk_min = np.argmax(dtest_val[0].mean(axis=1))
        idk_min_se = np.where(dtest_val[0].mean(axis=1) < target)
        done_se = np.min(dtest_list[idx:one_se])
        ddone_se = np.max(dtest_list[:one_se])
        ddone_se_r2 = np.min(dtest_list[one_se:r2])
        ddone_se_r2 = np.max(dtest_list[one_se:r2])

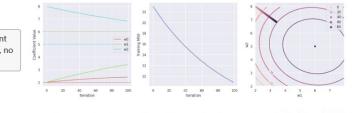
```



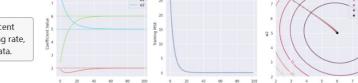
a. Batch gradient descent with large learning rate, no noise in training data.



b. Mini-batch (small batch size) gradient descent with moderate learning rate, high noise in training data.



c. Batch gradient descent with small learning rate, no noise in training data.



d. Batch gradient descent with moderate learning rate, no noise in training data.

W5 - Logistic regression and classification

Prediction: positive | Prediction: negative

Actual value positive	True positive $y=1, \hat{y}=1$	False negative $y=1, \hat{y}=0$
Actual value negative	False positive $y=0, \hat{y}=1$	True negative $y=0, \hat{y}=0$
Accuracy =	$\frac{TP+TN}{TP+FP+TN+FN} = \frac{TP+TN}{P+N}$	

$$P = TP + FN, N = FP + TN \quad (\text{actual})$$

Predicted-Positive	Actual - Positive	Actual - Negative
$PPV = \frac{TP}{TP+FP}$ (precision)	$FDR = \frac{FP}{TP+FP}$	
$TPR = \frac{TP}{TP+FN}$ (recall, sensitivity)	$FPR = \frac{FP}{FP+TN}$	

$$F_1 = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{balanced accuracy} = \frac{1}{2} \left(\frac{TP}{P} + \frac{TN}{N} \right)$$

email: spam, not spam

threshold $\hat{y} \rightarrow \hat{P}(y=1) \rightarrow \text{precision} \uparrow$
 $\downarrow \hat{P}(y=1) \text{ or stay} \rightarrow \text{recall for stay}$

$$\text{logistic/sigmoid function: } \sigma(z) = \frac{1}{1+e^{-z}}$$

$$z = \sum_{j=1}^n w_j x_{j,k}, P(y=1|x) = \frac{1}{1+e^{-z}}, P(y=0|x) = \frac{e^{-z}}{1+e^{-z}}$$

$$\text{softmax function: } g_k(z) = \frac{e^{z_k}}{\sum_{i=1}^n e^{z_i}} = P(y=k|x)$$

$$\text{Maximum likelihood estimate}$$

$$\hat{w} = \arg \max_w P(y|w) = P(y|x, w) = \prod_{i=1}^n P(y_i|w_i, w)$$

$$P(y_i|x_i, w) = y_i P(y_i=1|x_i, w) + (1-y_i) P(y_i=0|x_i, w)$$

Binary cross entropy loss function:

$$L(w) = - \sum_{i=1}^n (y_i \ln P(y_i=1|x_i, w) + (1-y_i) \ln P(y_i=0|x_i, w))$$

$$\frac{\partial L(w)}{\partial z} = \frac{1}{1+e^{-z}} - y_i, w^{fit} = w^0 + \sum_{i=1}^n (y_i - \ln w_i^0 / \sum_{j=1}^n w_j^0) x_i;$$

$$\text{Bayes rule: } P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

fitted logistic regression for binary classification

$$w = [-1.8, 7, -4, 4], \text{ first is bias}$$

$$2 \cdot \text{exp} = [0, 0, 0, 0, -0.02], z = w \cdot w + \sum_{i=1}^n w_i x_i = -1.8$$

$$P(y=1|x) = \frac{1}{1+e^{-z}} \approx 0.142$$

Scalar

```
# Scalar
scaler = StandardScaler().fit(Xtr)
Xtr_std = scaler.transform(Xtr)
Xts_std = scaler.transform(Xts)

vec = CountVectorizer(stop_words='english')
Xtr_vec = vec.fit_transform(Xtr_str)
Xts_vec = vec.transform(Xts_str)
```

Linear Regression

```
# Linear Regression
features = ['Park_Nearby', 'Grocery_Stores_Nearby', 'Schools_Nearby', 'Public_Transit_Nearby']
target = ['Availability_Score']
X = df[features]
y = df[target]
random_state = 23
Xtr, Xts, ytr, yts = train_test_split(X, y, test_size=0.2, random_state=random_state)
model = LinearRegression().fit(Xtr, ytr)
yts_hat = model.predict(Xts)
r2_sq = r2_score(yts, yts_hat)
mse = mean_squared_error(yts, yts_hat)

# Single split - random shuffle
random_state = 8
Xtr_one_shuf, Xts_one_shuf, ytr_one_shuf, yts_one_shuf = train_test_split(X, y, test_size=1/5, random_state=random_state)
model = LinearRegression()
model.fit(Xtr_one_shuf, ytr_one_shuf)
yts_one_shuf_pred = model.predict(Xts_one_shuf)
r2_kf_shuf = r2_score(yts_one_shuf, yts_one_shuf_pred)

# Single split - sorted data, no shuffle
Xtr_one_order, Xts_one_order, ytr_one_order, yts_one_order = train_test_split(X, y, test_size=1/5, shuffle=False)
model = LinearRegression()
model.fit(Xtr_one_order, ytr_one_order)
yts_one_order_pred = model.predict(Xts_one_order)
r2_one_order = r2_score(yts_one_order, yts_one_order_pred)

# Multiple splits - random shuffle
n_folds = 5
r2_kf_fold = np.zeros(shape=(n_fold,))
kf = KFold(n_splits=n_fold, shuffle=True, random_state=random_state)
for i, (idx_tr, idx_ts) in enumerate(kf.split(X)):
    model = LinearRegression()
    model.fit(X[idx_tr], ytr[idx_tr])
    y_pred_kf_fold = model.predict(X[idx_ts])
    r2_kf_fold[i] = r2_score(yidx_ts, y_pred_kf_fold)
    # Shuffle mean = np.mean(r2_kf_fold)
    # Multiple splits - time series
n_folds = 5
r2_ts = np.zeros(shape=(n_fold,))
ts = TimeSeriesSplit(n_splits=n_fold)
for i, (idx_tr, idx_ts) in enumerate(ts.split(X)):
    model = LinearRegression()
    model.fit(X[idx_tr], ytr[idx_tr])
    y_pred_ts = model.predict(X[idx_ts])
    r2_ts[i] = r2_score(yidx_ts, y_pred)
r2_ts_mean = np.mean(r2_ts)
```

Ridge Regression

```
x_names = ['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'x', 'y', 'z']
y_names = ['price']
random_state = 13
Xtr_df, Xts_df = train_test_split(df[x_names], test_size=0.3, random_state=random_state, shuffle=True)
ytr_df, yts_df = train_test_split(df[y_names], test_size=0.3, random_state=random_state, shuffle=True)
Xtr, Xts, ytr, yts = np.array(Xtr_df), np.array(Xts_df), np.array(ytr_df), np.array(yts_df)
alpha_list = np.array([10, 20, 50, 100, 200, 500])
nfold = 5
mse_val = np.zeros(len(alpha_list), nfold)
# k-fold
kf = KFold(n_splits=n_fold, shuffle=False)
# For each fold, standardize the data
for ifold, (idx_tr, idx_val) in enumerate(kf.split(Xtr)):
    X_train_fold_std = Xtr[idx_tr].copy(), Xtr[idx_val].copy()
    y_train_fold = ytr[idx_tr].copy(), ytr[idx_val].copy()
    X_val_fold_std = Xtr[idx_val].copy()
    X_train_fold_std = StandardScaler().fit(X_train_fold)
    X_train_fold_std = scaler.transform(X_train_fold)
    X_val_fold_std = scaler.transform(X_val_fold)
    # For each alpha in the list, fit a Ridge regression model on the standardized data
    for i, alpha in enumerate(alpha_list):
        model = Ridge(alpha=alpha)
        model.fit(X_train_fold_std, y_train_fold)
        y_pred = model.predict(X_val_fold_std)
        mse_val[i, ifold] = mean_squared_error(y_val_fold, y_pred)
mse_mean = np.mean(mse_val, axis=1)
alpha_min_mse = alpha_list[np.argmin(mse_mean)]
# entire training set
scaler = StandardScaler().fit(Xtr)
Xtr_std = scaler.transform(Xtr)
Xts_std = scaler.transform(Xts)
model = Ridge(alpha=alpha_min_mse)
model.fit(Xtr_std, ytr)
y_pred = model.predict(Xts_std)
mse_ridge = mean_squared_error(yts, y_pred)
```

K-fold CV with Fourier basis expansion

You decide to use a linear regression model with Fourier basis transformation of the `'hourofweek'` feature.

$$\hat{y} = w_0 + w_1x + \sum_{k \in \text{list}} w_{kx} \cos(2\pi x/k) + w_{kx} \sin(2\pi x/k)$$

where each sine and cosine pair represents the periodic behavior over a particular time interval.

For example, if `tlist = [0.5, 1]`, then your model would be:

$$\hat{y} = w_0 + w_1x + w_{0.5x} \cos(2\pi x/0.5) + w_{0.5x} \sin(2\pi x/0.5) + w_{1x} \cos(2\pi x/1) + w_{1x} \sin(2\pi x/1)$$

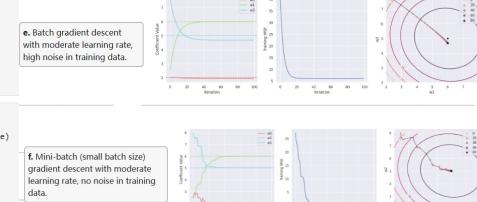
Now, you are ready to fit a K-fold CV in your CV, you will fit and evaluate a `LinearRegression` model (using the `sklearn` implementation) on an increasing number of columns of the data, as described above -

- In the first iteration of your K-fold CV, you will evaluate the regression for `tlist_eval = []`.
- In the second iteration of your K-fold CV, you will evaluate the regression for `tlist_eval = [1]`.
- In the third iteration, you will evaluate the regression for `tlist_eval = [1, 2]`.

and so on, until, in the final iteration, you will evaluate the regression for all the values in `tlist`.

(Of course, you won't re-compute the Fourier basis transformation inside the loop - you'll just select the appropriate rows and columns from `Xtr_trans` in each iteration.)

Since you have prepared a 'ones column' in the data, you will pass `fit_intercept=False` as an argument to the `LinearRegression`, so that it won't also fit another 'intercept' term (in addition to the coefficient for the 'ones' column).



```
# Linear Regression
df_tr, df_ts = train_test_split(df, train_size=10000, shuffle=False, random_state=42)
features = ['hour', 'month', 'dayofweek', 'hourofweek']
target = ['ridership']
model = LinearRegression()
model.fit(df_tr[features], df_tr[target])
y_pred = model.predict(df_ts[features])
r2_lr = r2_score(df_ts[target], y_pred)

# K-fold CV with Fourier basis expansion
tlist = []
Xtr = df_tr['hourofweek'].values
Xts = df_ts['hourofweek'].values
ytr = df_tr['ridership'].values
yts = df_ts['ridership'].values
Xtr_trans = np.column_stack((np.ones_like(Xtr), Xtr))
Xts_trans = np.column_stack((np.ones_like(Xts), Xts))

for t in tlist:
    cos_tr = np.cos(2 * np.pi * Xtr / t)
    sin_tr = np.sin(2 * np.pi * Xtr / t)
    cos_ts = np.cos(2 * np.pi * Xts / t)
    sin_ts = np.sin(2 * np.pi * Xts / t)
    Xtr_trans = np.column_stack((Xtr_trans, cos_tr, sin_tr))
    Xts_trans = np.column_stack((Xts_trans, cos_ts, sin_ts))

Xtr_trans = np.round(Xtr_trans, 10)
Xts_trans = np.round(Xts_trans, 10)

nfold = 5
r2_val = np.zeros((len(tlist) + 1, nfold))
kf = KFold(n_splits=n_fold, shuffle=False)
for i, t in enumerate(range(len(tlist) + 1)):
    # The first model uses only the intercept and original X (2 columns)
    # Subsequent models add Fourier features in pairs (cos and sin)
    num_columns = 2 + 2 * i
    X_subset = Xtr_trans[:, :num_columns]
    X_subset = Xts_trans[:, :num_columns]
    for fold, (train_idx, val_idx) in enumerate(kf.split(X_subset)):
        X_train, X_val = X_subset[train_idx], X_subset[val_idx]
        y_train, y_val = ytr[train_idx], ytr[val_idx]
        model = LinearRegression(fit_intercept=False)
        model.fit(X_train, y_train)
        y_pred = model.predict(X_val)
        r2_val[i, fold] = r2_score(y_val, y_pred)

r2_mean = r2_val.mean(axis=1)
r2_se = np.std(r2_val, axis=1, ddof=1) / np.sqrt(nfold)

idx_max = np.argmax(r2_mean)
# Compute the threshold for the one-SE rule
one_se = r2_mean[idx_max] - r2_se[idx_max]
# Find the simplest model (smallest number of Fourier features) whose R^2 is within one SE of the best model.
tlist_opt = []
for i in range(len(r2_mean)):
    if r2_mean[i] > one_se:
        tlist_opt.append(tlist[i]) # Select the first i values of tlist
        break
# Train the model on the entire training set for this 'tlist_opt', then evaluate its performance on the test set.
# Save the R^2 score in r2_one_se.
tlist_opt.append(2 + 2 * len(tlist_opt))
Xtr_opt = Xtr_trans[:, :num_columns]
Xts_opt = Xts_trans[:, :num_columns]
model_opt = LinearRegression(fit_intercept=False)
model_opt.fit(Xtr_opt, ytr)
y_pred_opt = model_opt.predict(Xts_opt)
r2_one_se = r2_score(yts, y_pred_opt)
```

Logistic Regression

```
# Standardize data in each fold.
random_state = 14
Xtr, Xts, ytr, yts = train_test_split(X, y, test_size=0.3, random_state=random_state)
med_values = np.nanmedian(Xtr, axis=0)
Xtr.fillna(med_values, inplace=True)
Xts.fillna(med_values, inplace=True)
C_best = np.logspace(-1,3,10)
nFold = 3
acc_mean = np.zeros((len(C_best), nFold))
for iC, C in enumerate(C_best):
    kf = KFold(n=nFold, shuffle=False)
    for iFold, (Itr, Ival) in enumerate(kf.split(Xtr)):
        Xtr_fold, Xval_fold = Xtr.iloc[Itr], Xtr.iloc[Ival]
        ytr_fold, yval_fold = ytr[Itr], ytr[Ival]
        Xtr_std = StandardScaler().fit(Xtr_fold)
        Xtr_fold = Xtr_std.transform(Xtr_fold)
        Xval_std = StandardScaler().fit(Xval_fold)
        Xval_fold = Xval_std.transform(Xval_fold)
        clf = LogisticRegression(random_state = random_state, solver = 'liblinear', penalty='l1', C = C)
        clf.fit(Xtr_fold, ytr_fold)
        yhat = clf.predict(Xval_fold)
        acc_val[iC, ifold] = accuracy_score(yval_fold, yhat)
    acc_mean[iC, :] = np.mean(acc_val[:, axis=1])
C_best = C[test(np.argmax(acc_mean))]
# Create a vector of words in English
vec = CountVectorizer(stop_words='english')
Xtr_vec = vec.fit_transform(Xtr)
Xts_vec = vec.transform(Xts)
C_best = np.logspace(-3, 3, num=20)
nFold = 5
acc_mean = np.zeros((len(C_best), nFold))
for iC, C in enumerate(C_best):
    kf = KFold(n=nFold, shuffle=False)
    for iFold, (Itr, Ival) in enumerate(kf.split(Xtr_vec)):
        for IC, C in enumerate(vec.get_feature_names()):
            clf = LogisticRegression(random_state = 0, penalty='l1', solver='liblinear', C = C)
            clf.fit(Xtr_vec, ytr)
            yhat = clf.predict(Xval_vec)
            acc_val[IC, ifold] = accuracy_score(ytr[Ival], yhat)
        acc_mean = np.mean(acc_val[:, axis=1])
C_best = C[test(np.argmax(acc_mean))]
model_best = LogisticRegression(penalty='l1', C=C_best, random_state=0, solver='liblinear')
model_best.fit(Xtr_vec, ytr)
y_pred_best = model_best.predict(Xts_vec)
acc_best = accuracy_score(yts, y_pred_best)
count_best = np.count_nonzero(model_best.coef_)
acc_mean = np.mean(acc_mean[:, axis=1])
C.one_se = C * test(np.sqrt(np.log(len(C)*2)))
model_one_se = LogisticRegression(penalty='l1', C=C.one_se, random_state=0, solver='liblinear')
model_one_se.fit(Xtr_vec, ytr)
y_pred_one_se = model_one_se.predict(Xts_vec)
acc_one_se = accuracy_score(yts, y_pred_one_se)
count_one_se = np.count_nonzero(model_one_se.coef_)
```

Predicting and Preventing Deaths in the ICU

```
]: # Proportion
X = df[col_names[-1]]
y = df['In-hospital_death'].values
Xtr, Xts, ytr, yts = train_test_split(X, y, test_size=800, shuffle=True, random_state=0)
Xtr, Xts, ytr, yts = train_test_split(Xtr, ytr, test_size=800, shuffle=True, random_state=0)
col_meds = Xtr.median()
Xtr = Xtr.fillna(col_meds)
Xts = Xts.fillna(col_meds)
y = np.sum(yts)
clf = LogisticRegression(random_state = 0, penalty=None)
clf.fit(Xtr, ytr)
y_pred_lr = clf.predict_proba(Xts)
acc_lr = accuracy_score(yts, y_pred_lr)
most_frequent_label = np.bincount(ytr).argmax()
y_pred_base = np.full_like(yts, most_frequent_label)
acc_base = accuracy_score(yts, y_pred_base)
# ROC curve
yprob = clf.predict_proba(Xts)[:, 1]
fpr, tpr, thresholds = roc_curve(yts, yprob)
plt.plot(fpr, tpr)
plt.xlabel('False Positive Rate')
plt.xlabel('True Positive Rate')
plt.title('ROC Curve')
plt.show()
```

Evaluate on the test set

```
mort_ts = np.sum(yts)
# Policy 0: no review, so no lives are saved
saved_0 = 0
# Policy 1: review all
N = len(yts)
total_time = (N / 25) * 60
t_ts_1 = total_time / N
p_ts_1 = np.exp(-(t_ts_1 ** 2) / 100)
at_risk_patients = np.sum(yts)
mort_ts_1 = at_risk_patients * p_ts_1
saved_1 = mort_ts - mort_ts_1
# Policy 2: review using 50% threshold
y_prob_ts = clf.predict_proba(Xts)[:, 1]
high_risk_patients = y_prob_ts > 0.5
num_high_risk = np.sum(high_risk_patients)
t_ts_2 = total_time / num_high_risk if num_high_risk > 0 else 0
p_ts_2 = np.exp(-(t_ts_2 ** 2) / 100)
at_risk_high_risk = np.sum(yts[high_risk_patients])
at_risk_low_risk = np.sum(yts[~high_risk_patients])
mort_ts_2 = at_risk_high_risk * p_ts_2 + at_risk_low_risk
saved_2 = mort_ts - mort_ts_2
# Policy 3: review using the opt
high_risk_patients = y_prob_ts > thr_opt
num_high_risk = np.sum(high_risk_patients)
t_ts_3 = total_time / num_high_risk if num_high_risk > 0 else 0
p_ts_3 = np.exp(-(t_ts_3 ** 2) / 100)
at_risk_high_risk = np.sum(yts[high_risk_patients])
at_risk_low_risk = np.sum(yts[~high_risk_patients])
mort_ts_3 = at_risk_high_risk * p_ts_3 + at_risk_low_risk
saved_3 = mort_ts - mort_ts_3
```

If the doctor spends t minutes reviewing a file, the probability of missing the correct treatment is $P(\text{fail}) = \exp(-t^2/100)$.

Policy 0: No extra review

Policy 1: Review all files

Now suppose the specialist team will spend

$$\frac{N}{25} \times 60$$

minutes total on extra review, where N is the number of files in the validation set, and they will spend an equal amount of time on each of the N files in the validation set.

Policy 2: Review high-risk files with 0.5 threshold

Now, suppose that the doctor only reviews the files of patients for whom the model outputs a probability of mortality greater than 50%. They will spend the same total amount of time as before,

$$\frac{N}{25} \times 60$$

Compute:

- \bar{x}_2 = the time (in minutes) spent on each file
- p_2 = the probability of failing to identify the correct intervention, given this \bar{x} and then use them to find the expected number of deaths in the validation set with this policy. See this value in `mort_2`.

Policy 3: Review high-risk files with optimal threshold

For all thresholds in `np.arange(0, 1, 0.01)`, compute the expected mortality, given the doctor only reviews the files of patients for whom the model outputs a probability of mortality greater than that threshold. They will spend the same total amount of time as before,

$$\frac{N}{25} \times 60$$

minutes total on extra review, where N is the number of files in the validation set. Save the results in `mort_thr`.

Policy 0: No extra review
mort_0 = np.sum(yts)

Policy 1: Review all files
N = len(yts)
Step 1: Compute time spent per file
t_1 = (N / 25) * 60 # Minutes per file
Step 2: Compute probability of missing correct intervention
p_1 = np.exp(-t_1**2 / 100)
Step 3: Compute expected mortality under Policy 1
mort_1 = np.sum(yts * p_1)
Step 4: Predict probabilities of mortality for the validation set

Policy 2: Review high-risk files with 0.5 threshold
y_prob_a1 = clf.predict_proba(Xts)[:, 1]
Identify high-risk and low-risk patients
high_risk_indices = [i for i in range(len(yts)) if y_prob_a1[i] > 0.5]
low_risk_indices = [i for i in range(len(yts)) if y_prob_a1[i] <= 0.5]
Number of high-risk files
N_high_risk = len(high_risk_indices)
Step 2: Calculate t_2
N = len(yts)
t_2 = (N / 25) * 60
t_2 = T / N_high_risk
Step 3: Calculate p_2
p_2 = np.exp(-t_2**2 / 100)
Step 4: Estimate expected deaths
mort_2 = 0
for i in high_risk_indices:
 if yts[i] == 1: # Patient would have died without review
 mort_2 += p_2 # Adjusted probability of death
 else:
 mort_2 += 0 # No death for survivors
for i in low_risk_indices:
 mort_2 += yts[i] # No review, so deaths remain unchanged

Policy 3: Review high-risk files with optimal threshold
Step 1: Predict probabilities of mortality for the validation set
y_prob_a1 = clf.predict_proba(Xts)[:, 1] # Predicted probabilities for validation set
Total number of files in the validation set
N = len(yts)
Total time spent on extra review (same as in Policy 1 and Policy 2)
T = (N / 25) * 60
Step 2: Initialize arrays to store results
thr_list = np.arange(0, 1, 0.01) # List of thresholds

Find

- `thr_opt`, the threshold that minimizes mortality on the validation set, and
- `mort_3`, the expected mortality on the validation set with this threshold.

```
thr_opt = thr_list[np.argmin(mort_thr)]  
# mort_3 = np.min(mort_thr)  
mort_3 = np.min(mort_thr[mort_thr > 0]) if np.any(mort_thr > 0) else np.min(mort_thr)
```