

- A supervised machine learning "recipe"
- Step 1: Get labeled data: $(x_i, y_i), i = 1, 2, \dots, n$.
 - Step 2: Choose a candidate model f : $\hat{y} = f(x)$.
 - Step 3: Select a loss function.
 - Step 4: Find the model parameter values that minimize the loss function (training).
 - Step 5: Use trained model to predict \hat{y}' for new samples not used in training (inference).
 - Step 6: Evaluate how well your model generalizes.

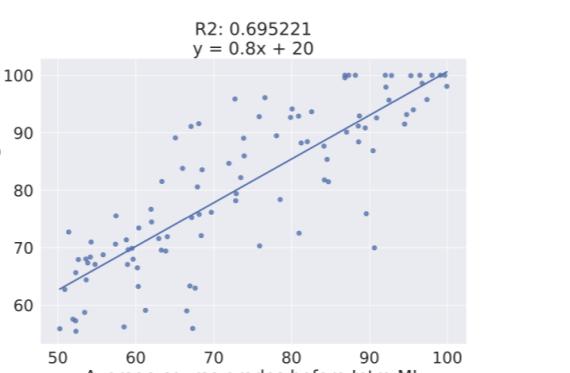


Figure 5: Predicting students' grades in Intro ML using regression on previous coursework.

To some extent, a student's average grades on previous coursework "explains" their grade in Intro ML.

- The predicted value for each student, \hat{y} , is along the diagonal line. Draw a vertical line from each student's point (y) to the corresponding point on the line (\hat{y}). This is the residual $e = y - \hat{y}$.
- Some students fall right on the line - these are examples that are explained "well" by the model.
- Some students are far from the line. The magnitude of the *residual* is greater for these examples.
- The difference between the "true" value y and the predicted value \hat{y} may be due to all kinds of differences between the "well-explained example" and the "not-well-explained-example" - not everything about Intro ML course grade can be explained by performance in previous coursework! This is what the residual captures.

Interpreting the linear regression: If slope w_1 is 0.8 points in Intro ML per point average in previous coursework, we can say that

- a 1-point increase in score on previous coursework is, on average, associated with a 0.8 point increase in Intro ML course grade.

What can we say about possible explanations? We can't say much using this method - anything is possible:

- statistical fluke (we haven't done any test for significance)
- causal - students who did well in previous coursework are better prepared
- confounding variable - students who did well in previous coursework might have more time to study because they don't have any other jobs or obligations, and they are likely to do well in Intro ML for the same reason.

This method doesn't tell us *why* this association is observed, only that it is. (There are other methods in statistics for determining whether it is a statistical fluke, or for determining whether it is a causal relationship.)(Also note that the 0.8 point increase according to the regression model is only an *estimate* of the "true" relationship.)

Linear basis function model for regression

Standard linear model:

$$\hat{y}_i = w_0 + w_1 x_{i,1} + \dots + w_d x_{i,d}$$

Linear basis function model:

$$\hat{y}_i = w_0 \phi_0(\mathbf{x}_i) + \dots + w_p \phi_p(\mathbf{x}_i)$$

Some notes:

- The 1s column we added to the design matrix is easily represented as a basis function ($\phi_0(\mathbf{x}) = 1$).
- There is not necessarily a one-to-one correspondence between the columns of X and the basis functions ($p \neq d$ is OK). You can have more/fewer basis functions than columns of X .
- Each basis function can accept as input the entire vector \mathbf{x}_i .
- The model has $p+1$ parameters.

Vector form of linear basis function model

The prediction of this model expressed in vector form is:

$$\hat{y}_i = (\phi(\mathbf{x}_i), \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}_i)$$

where

$$\phi(\mathbf{x}_i) = [\phi_0(\mathbf{x}_i), \dots, \phi_p(\mathbf{x}_i)], \mathbf{w} = [w_0, \dots, w_p]$$

(The angle brackets denote a dot product.)

Important note: although the model can be non-linear in \mathbf{x} , it is still linear in the parameters \mathbf{w} (note that \mathbf{w} appears outside $\phi(\cdot)$). That's what makes it a *linear model*.

Some basis functions have their own parameters that appear inside the basis function, i.e. we might have a model

$$\hat{y}_i = \mathbf{w}^T \phi(\mathbf{x}_i, \theta)$$

where θ are the parameters of the basis function. The model is *non-linear* in those parameters, and they need to be fixed before training.

Matrix form of linear basis function model

Given data $(\mathbf{x}_i, y_i), i = 1, \dots, n$:

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_p(\mathbf{x}_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \dots & \phi_p(\mathbf{x}_n) \end{bmatrix}$$

and $\hat{\mathbf{y}} = \Phi \mathbf{w}$.

Interpreting R2: coefficient of determination (1)

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

For linear regression: What proportion of the variance in y is "explained" by our model?

- $R^2 \approx 1$ - model "explains" all the variance in y
- $R^2 \approx 0$ - model doesn't "explain" any of the variance in y

Interpreting R2: coefficient of determination (2)

Alternatively: what is the ratio of error of our model, to error of prediction by mean?

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

Gradient of MSE

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2$$

gives us the gradient

$$\nabla L(\mathbf{w}) = -\Phi^T (\mathbf{y} - \Phi \mathbf{w})$$

Solving for \mathbf{w}

$$\begin{aligned} \nabla L(\mathbf{w}) &= 0, \\ -\Phi^T (\mathbf{y} - \Phi \mathbf{w}) &= 0, \\ \Phi^T \Phi \mathbf{w} &= \Phi^T \mathbf{y}, \text{ or} \\ \mathbf{w} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \end{aligned}$$

Solving a set of linear equations

If $\Phi^T \Phi$ is full rank (usually: if $n \geq d$), then a unique solution is given by

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

This expression:

$$\Phi^T \mathbf{w} = \Phi^T \mathbf{y}$$

represents a set of d equations in d unknowns, called the *normal equations*.

We can solve this as we would any set of linear equations (see supplementary notebook on computing regression coefficients by hand).

Week 2 Labs

Compare this to the known solution:
 $1 \text{ np.linalg.inv(A.T.dot(A))}.dot(A.T.dot(y))$
 $\text{array([5.58391608, 0.77972028, -1.6993007])$

The optimal least squares values for the vector \mathbf{w} are
 $\mathbf{w}^* = (A^T A)^{-1} A^T \mathbf{y}$ Note that the least-squares solutions are the solutions of the matrix equation
 $A^T A \mathbf{w} = A^T \mathbf{y}$

Note: other ways to do the same thing...

```
# first, add a ones column to design matrix
x_tilde = np.hstack((np.ones(n_samples, 1), x_train))
# using matrix operations to find w = (X^T X)^{-1} X^T y
print( (np.linalg.inv(x_tilde.T @ x_tilde) @ x_tilde.T @ y_train) )
```

```
# using solve on normal equations: X^T X w = X^T y
# solve only works on matrix that is square and of full-rank
# see https://numpy.org/doc/stable/reference/generated/numpy.linalg.solve.html
print( np.linalg.solve(x_tilde.T @ x_tilde, x_tilde.T @ y_train) )
```

```
# using the lstsq solver
# problem may be under-, well-, or over-determined
# see https://numpy.org/doc/stable/reference/generated/numpy.linalg.lstsq.html
print( np.linalg.lstsq(x_tilde, y_train, rcond=0)[0] ) xt
```

Week 2 Linear Regression

Example: Intro ML grades (2)

Given

$$MSE(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

we take

$$\frac{\partial MSE}{\partial w_0} = 0, \quad \frac{\partial MSE}{\partial w_1} = 0$$

Optimizing w - simple linear regression (2)

First, the intercept:

$$MSE(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial MSE}{\partial w_0} = -2 \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))$$

using chain rule, power rule.

(We can then drop the -2 constant factor when we set this expression equal to 0.)

Optimizing w - simple linear regression (3)

Set this equal to 0, "distribute" the sum, and we can see

$$\frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) = 0$$

$$\Rightarrow w_0^* = \bar{y} - w_1^* \bar{x}$$

where \bar{x}, \bar{y} are the means of x, y .

Optimizing w - simple linear regression (4)

Now, the slope coefficient:

$$MSE(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial MSE}{\partial w_1} = \frac{1}{n} \sum_{i=1}^n 2(y_i - w_0 - w_1 x_i)(-x_i)$$

11

Optimizing w - simple linear regression (5)

$$\Rightarrow -2 \frac{1}{n} \sum_{i=1}^n x_i (y_i - w_0 - w_1 x_i) = 0$$

Solve for w_1^* :

$$w_1^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Note: some algebra is omitted here, but refer to the secondary notes for details.

Optimizing w - simple linear regression (6)

The slope coefficient is the ratio of sample covariance σ_{xy} to sample variance σ_x^2 :

$$\frac{\sigma_{xy}}{\sigma_x^2}$$

where $\sigma_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ and $\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Optimizing w - simple linear regression (7)

We can also express it as

$$\frac{r_{xy} \sigma_y}{\sigma_x}$$

where sample correlation coefficient $r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$.(Note: from Cauchy-Schwartz law, $|\sigma_{xy}| < \sigma_x \sigma_y$, we know $r_{xy} \in [-1, 1]$)

MSE for optimal simple linear regression

$$MSE(w_0^*, w_1^*) = \sigma_y^2 - \frac{\sigma_{xy}^2}{\sigma_x^2}$$

$$\Rightarrow \frac{MSE(w_0^*, w_1^*)}{\sigma_y^2} = 1 - \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2}$$

- the ratio on the left is the *fraction of unexplained variance*: of all the variance in y , how much is still "left" unexplained after our model explains some of it? (best case: 0)
- the ratio on the right is the *coefficient of determination*, R^2 (best case: 1).

Ordinary least squares solution for multiple/linear basis function regression

Setup: L2 norm

Definition: L2 norm of a vector $\mathbf{x} = (x_1, \dots, x_n)$:

$$\|\mathbf{x}\| = \sqrt{x_1^2 + \dots + x_n^2}$$

We will want to minimize the L2 norm of the residual.

Setup: Gradient vector

To minimize a multivariate function $f(\mathbf{x}) = f(x_1, \dots, x_n)$, we find places where the *gradient* is zero, i.e. each entry must be zero:For linear regression: What proportion of the variance in y is "explained" by our model?

- $R^2 \approx 1$ - model "explains" all the variance in y
- $R^2 \approx 0$ - model doesn't "explain" any of the variance in y

Interpreting R2: coefficient of determination (2)

Alternatively: what is the ratio of error of our model, to error of prediction by mean?

MSE for multiple/LBF regression

Given a vector \mathbf{y} and matrix Φ (with d columns, n rows),

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2$$

where the norm above is the L2 norm.

Week 4 Model Selection

- Why use a subset of features?
 - High risk of overfitting if you use all features!
 - For linear regression, there's a unique OLS solution only if $n \geq d$
 - For linear regression, when $N \geq p$, variance increases linearly with number of parameters, inversely with number of samples.
- Problems with hold-out validation
- Fitted model (and test error) varies a lot depending on samples selected for training and validation.
 - Fewer samples available for estimating parameters.
 - Especially bad for problems with small number of samples.

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

This is not like K-fold CV which uses non-overlapping validation sets (they are only the same for $p = 1$)!

Computation (leave-p-out CV)

$\binom{n}{p}$ iterations, in each:

- train on $n - p$ samples
- score on p samples

Usually, this is too expensive - but sometimes LOO CV can be a good match to the model (KNN).

Computation (K-fold CV)

K iterations, in each:

- train on $n - n/k$ samples
- score on n/k samples

K-fold CV - how to split?

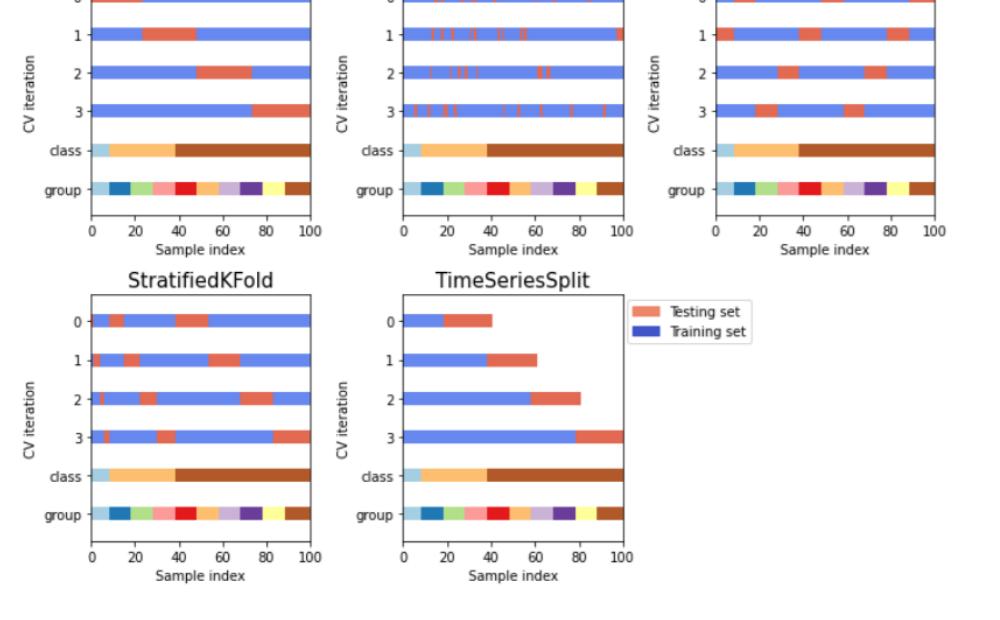


Figure 6: K-fold CV variations.

Selecting the right K-fold CV is very important for avoiding data leakage! (Also for training/test split.)

Refer to the [function documentation](#) for more examples.

One standard error rule

- Model selection that minimizes mean error often results in too-complex model
- One standard error rule: use simplest model where mean error is within one SE of the minimum mean error

One standard error rule - algorithm (1)

- Given data X, y
- Compute score $S_{p,i}$ for model p on fold i (of K)
- Compute average (\bar{S}_p), standard deviation σ_p , and standard error of scores:

$$SE_p = \frac{\sigma_p}{\sqrt{K-1}}$$

One standard error rule - algorithm (2)

"Best score" model selection: $p^* = \arg\min_p \bar{S}_p$

One SE rule for "lower is better" scoring metric: Compute target score: $S_t = \bar{S}_{p^*} + SE_{p^*}$

then select simplest model with score lower than target:

$$p^{*,\text{SE}} = \min\{p | \bar{S}_p \leq S_t\}$$

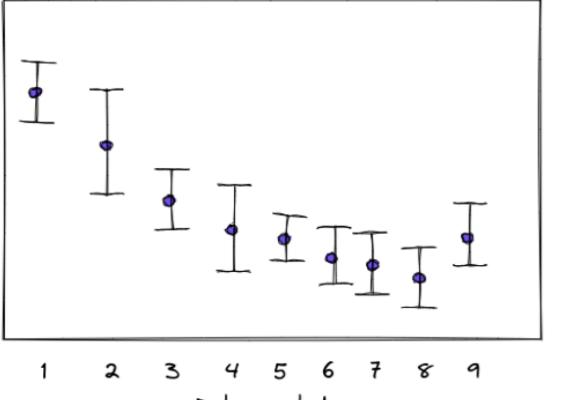


Figure 12: Model selection using one SE rule on MSE. The best scoring model is $d = 8$, but $d = 6$ is simplest model within one SE of the best scoring model, and so $d = 6$ would be selected according to the one-SE rule.

Note: this assumes you are using a "smaller is better" metric such as MSE. If you are using a "larger is better" metric, like R2, how would we change the algorithm?

One standard error rule - algorithm (3)

"Best score" model selection: $p^* = \arg\max_p \bar{S}_p$

One SE rule for "higher is better" scoring metric: Compute target score: $S_t = \bar{S}_{p^*} - SE_{p^*}$

then select simplest model with score higher than target:

$$p^{*,\text{SE}} = \min\{p | \bar{S}_p \geq S_t\}$$

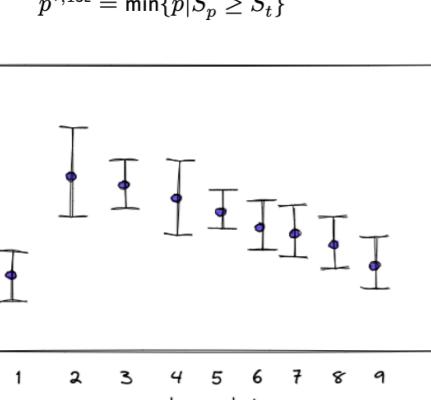


Figure 13: Model selection using one SE rule on R2. In this example, the best scoring model is $d = 2$, and there is no simpler model within one SE, so the one-SE rule would also select $d = 2$.

Regularization

Penalty for model complexity

With no bounds on complexity of model, we can always get a model with zero training error on finite training set - overfitting.

Basic idea: apply penalty in loss function to discourage more complex models

Regularization vs. standard LS

Least squares estimate:

$$\hat{w} = \underset{w}{\operatorname{arg\min}} MSE(w), \quad MSE(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Regularized estimate w/ regularizing function $\phi(w)$:

$$\hat{w} = \underset{w}{\operatorname{arg\min}} J(w), \quad J(w) = MSE(w) + \phi(w)$$

Common regularizers

Ridge regression (L2):

$$\phi(w) = \alpha \sum_{j=1}^d |w_j|^2$$

LASSO regression (L1):

$$\phi(w) = \alpha \sum_{j=1}^d |w_j|$$

Common features: Ridge and LASSO

- Both penalize large w_j
- Both have parameter α that controls level of regularization
- Intercept w_0 not included in regularization sum (starts at 1!), this depends on mean of y and should not be constrained.

Differences: Ridge and LASSO (1)

- Ridge (L2):
 - minimizes $|w|^2$,
 - minimal penalty for small non-zero coefficients
 - heavily penalizes large coefficients
 - tends to make many "small" coefficients • Not for feature selection

Differences: LASSO (2)

- LASSO (L1):
 - minimizes $|w|$
 - tends to make coefficients either 0 or large (sparse!)
 - does feature selection (setting w_j to zero is equivalent to un-selecting feature)

Standardization (1)

Before learning a model with regularization, we typically standardize each feature and target to have zero mean, unit variance:

$$\begin{aligned} x_{i,j} &\rightarrow \frac{x_{i,j} - \bar{x}_j}{s_{x_j}} \\ y_i &\rightarrow \frac{y_i - \bar{y}}{s_y} \end{aligned}$$

Standardization (2)

Why?

- Without scaling, regularization depends on data range
- With mean removal, no longer need w_0 , so regularization term is just L1 or L2 norm of coefficient vector

L1 and L2 norm with standardization (1)

Assuming data standardized to zero mean, unit variance, the Ridge cost function is:

$$\begin{aligned} J(w) &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^d |w_j|^2 \\ &= ||\mathbf{A}w - \mathbf{y}||^2 + \alpha ||w||^2 \end{aligned}$$

L1 and L2 norm with standardization (2)

LASSO cost function ($||w||_1$ is L1 norm):

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

= $||\mathbf{A}w - \mathbf{y}||^2 + \alpha ||w||_1$

Ridge regularization

Why minimize $||w||^2$?

Without regularization:

- large coefficients lead to high variance
- large positive and negative coefficients cancel each other for correlated features (remember attractiveness ratings in linear regression case study...)

Ridge term and derivative

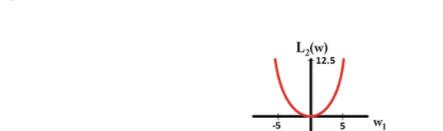


Figure 2: L2 term and its derivative for one parameter

Ridge closed-form solution

Taking derivative:

$$\frac{\partial J(w)}{\partial w} = 2\mathbf{A}^T(\mathbf{y} - \mathbf{Aw}) + 2\alpha w$$

Setting it to zero, we find

$$\mathbf{w}_{\text{ridge}} = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$$

LASSO term and derivative

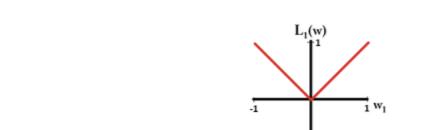


Figure 3: L1 term and its derivative for one parameter

Effect of regularization level

Greater α , more complex model:

- more α makes coefficients smaller
- LASSO: Greater α makes more weights zero.

Selecting regularization level

How to select α by CV:

- Outer loop: loop over CV folds
- Inner loop: loop over α

Delay Function

def create_dly_data(X,y,dly):

Create delayed data

...= X.shape

Xdly = np.zeros((n-dly,(dly+1)*p)) for i in range(dly+1):

Xdly[:,i:(i+dly)] = X[:,i:(i+dly)]

ydly = y[dly:]

return Xdly, ydly

Week 4 Model Selection

Week 4 Model SelectionNB

def generate_polynomial_regression_data(n=100, xrange=[-1,0,5,0,2], sigma=0.5):

x = np.random.uniform(xrange[0], xrange[1], n)

y = np.polynomial.polynomial.polyval(x,coefs) +

sigma * np.random.randn(n)

return x.reshape(-1,1), y

KFold using scikit vs for loop

kf = KFold(n_splits=nfold,shuffle=True)

for isplit, idx in enumerate(kf.split(x_train)):

idx_tr, idx_val = idx

this is same as

nfold = 5

number of folds (you choose)

nval = x_train.shape[0]

/nfold # number of validation samples per fold

idx_split = [i*nval for i in range(nfold)]

idx_list = np.arange(x_train.shape[0])

list of training data indices

np.random.shuffle(idx_list)

shuffle list of indices

for i, idx in enumerate(idx_split):

idx_val = idx_list[idx-1:nval]

idx_tr = np.delete(idx_list, idx_val)

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 how to split the data for isplit, idx in enumerate(idxval):

x_train_kfold = df_tr_poly[idx]

y_train_kfold = df_tr_deathincrease.values[idx]

x_val_kfold = df_tr_poly[idx:nval]

y_val_kfold = df_tr_deathincrease.values[idx:nval]

get transformed features

x_train_dtest = x_train_kfold[:, :dtest]

x_val_dtest = x_val_kfold[:, :dtest]

fit data

reg_dtest = LinearRegression().fit(x_train_dtest, y_train_kfold)

measure MSE on validation data

y_hat = reg_dtest.predict(x_val_dtest)

mse_val[idx, isplit] = metrics.mean_squared_error(y_val_kfold, y_hat)

r2_val[idx, isplit] = metrics.r2_score(y_val_kfold, y_hat)

measure MSE on training data

y_hat_tr = reg_dtest.predict(x_train_dtest)

Transformations, bias, variance

- Can use basis functions to map problem to transformed feature space (if "natural" decision boundary is non-linear)
- Variance increases with d and decreases with n
- Can add a regularization penalty to loss function

"Recipe" for logistic regression (binary classifier)

- Choose a model:

$$P(y=1|x, w) = \sigma\left(w_0 + \sum_{i=1}^d w_i x_i\right)$$
$$\hat{y} = \begin{cases} 1, & P(y=1|x) \geq t \\ 0, & P(y=1|x) < t \end{cases}$$

- Get data - for supervised learning, we need labeled examples: $(x_i, y_i), i = 1, 2, \dots, n$
- Choose a loss function that will measure how well model fits data: binary cross-entropy

$$\sum_{i=1}^n \ln(1 + e^{z_i}) - y_i z_i$$

- Find model parameters that minimize loss: use numerical optimization to find weight vector w
- Use model to predict \hat{y} for new, unlabeled samples.

"Recipe" for logistic regression (multi-class classifier)

- Choose a model: find probability of belonging to each class, then choose the class for which the probability is highest.

$$P(y=k|x) = \frac{e^{z_k}}{\sum_{\ell=1}^K e^{z_\ell}} \text{ where } \mathbf{z} = \mathbf{W}\mathbf{x}$$

- Get data - for supervised learning, we need labeled examples: $(x_i, y_i), i = 1, 2, \dots, n$
- Choose a loss function that will measure how well model fits data: cross-entropy

$$\sum_{i=1}^n \left[\ln \left(\sum_k e^{z_{ik}} \right) - \sum_k z_{ik} r_{ik} \right] \text{ where}$$

$$r_{ik} = \begin{cases} 1 & y_i = k \\ 0 & y_i \neq k \end{cases}$$

- Find model parameters that minimize loss: use numerical optimization to find weight vector w
- Use model to predict \hat{y} for new, unlabeled samples.

Generate data for logistic regression

```
# Generate training data
X = np.random.uniform(0, 1, size=(n_samples, 2))
y = np.array(w_true[0]+w_true[1]*X[:, 0]+w_true[2]*X[:, 1]>=0).astype(int)
```

Add some noise?

```
X = X + sigma*np.random.randn(X.shape[0])*X.shape[1].reshape(X.shape[0], X.shape[1])
```

Logistic regression training

```
def sigmoid(z):
    return 1/(1+np.exp(-z))

def cross_entropy_loss(w, X, y):
    p_pred = sigmoid(np.dot(X, w))
    loss_pos = y*np.log(p_pred)
    loss_neg = (1-y)*np.log(1-p_pred)
    return np.sum(-1*(loss_pos + loss_neg))

def gd_step(w, X, y, lr):
    p_pred = sigmoid(np.dot(X, w))
    gradient = np.dot(X.T, y - p_pred)
    w = w - lr * gradient
    l = cross_entropy_loss(w, X, y)
    return w, l, gradient
```

```
itr = 20000
lr = 0.005
tol = 0.01
w_init = np.random.randn(3)
print(w_init)
w_steps = np.zeros((itr, len(w_init)))
l_steps = np.zeros(itr)
grad_steps = np.zeros((itr, len(w_init)))
stop = 0
```

```
w_star = w_init
for i in range(itr):
    w_star, loss, grad = gd_step(w_star, X, y, lr)
    w_steps[i] = w_star
    l_steps[i] = loss
    grad_steps[i] = grad
    if np.linalg.norm(grad, ord=1) <= tol:
        stop = 1
        print("Stopping gradient descent at iteration %d" % i)
```

What can we do, as machine learning engineers, that might help?

Exploratory data analysis - look for possible underlying bias in the data.

Avoid using sensitive group as a feature (or a proxy for a sensitive group), but this doesn't necessarily help. (It didn't help in this example.)

Sometimes, using the sensitive group to explicitly add fairness might be better.

Make sure different groups are well represented in the data.

End users must be made aware of what the model output means -

for example, judges should understand that a "high risk" label only means a 65% chance of reoffending.

Evaluate final model for bias with respect to sensitive groups.

May not be possible to satisfy multiple fairness metrics simultaneously, work with end users to decide which fairness metrics to prioritize, and to create a model that is fair w.r.t. that metric.

xrange=[-1, 1]

```
def generate_polynomial_classifier_data(n=100, xrange=[-1, 1], coefs=[1, 0, 5, 0, 2], sigma=0.5):
    x = np.random.uniform(xrange[0], xrange[1], size=(n, 2))
    ysep = np.polynomial.polynomial.polyval(x[:, 0], coefs)
    y = (x[:, 1]>ysep).astype(int)
```

x[:, 0] = x[:, 0] + sigma * np.random.randn(n)
x[:, 1] = x[:, 1] + sigma * np.random.randn(n)

return x, y

clf = LogisticRegression(penalty='none', tol=0.01, solver='saga')

clf.fit(X, y)

clf.score(X, y)

Notes on logistic regression

- When there is under-modeling, adding transformed features can reduce bias.
- However, adding features also increases variance.
- We can use cross validation to choose a model that fits the data well but also generalizes.
- We can add a regularization penalty to our loss function to reduce variance.

Confusion matrix from df

```
cm = pd.crosstab(df['is_med_or_high_risk'], df['two_year_recid'],
                  rownames=['Predicted'], colnames=['Actual'])
p = plt.figure(figsize=(5, 5));
p = sns.heatmap(cm, annot=True, fmt="d", cbar=False)
```

Naive Bayes classifier

A quick look at a different type of model!

Probabilistic models (1)

For logistic regression, minimizing the cross-entropy loss finds the parameters for which $P(\mathbf{y}|\mathbf{X}, \mathbf{w})$ is maximized.

Probabilistic models (2)

For linear regression, assuming normally distributed stochastic error, minimizing the squared error loss finds the parameters for which $P(\mathbf{y}|\mathbf{X}, \mathbf{w})$ is maximized.

Surprise! We've been doing maximum likelihood estimation all along.

Probabilistic models (3)

ML models that try to

- get a good fit for $P(\mathbf{y}|\mathbf{X})$: discriminative models.
- fit $P(\mathbf{X}, \mathbf{y})$ or $P(\mathbf{X}|\mathbf{y})P(\mathbf{y})$: generative models.

Linear regression and logistic regression are both considered discriminative models; they say "given that we have this data, what's the most likely label?" (e.g. learning a mapping from an input to a target variable).

Generative models try to learn "what does data for each class look like" and then apply Bayes rule.

Bayes rule

For a sample \mathbf{x}_i , y_k is label of class k :

$$P(y_k|\mathbf{x}_i) = \frac{P(\mathbf{x}_i|y_k)P(y_k)}{P(\mathbf{x}_i)}$$

- $P(y_k|\mathbf{x}_i)$: posterior probability: "What is the probability that this sample belongs to class k , given its observed feature values are \mathbf{x}_i ?"
- $P(\mathbf{x}_i|y_k)$: conditional probability: "What is the probability of observing the feature values \mathbf{x}_i in a sample, given that the sample belongs to class k ?"
- $P(y_k)$: prior probability
- $P(\mathbf{x}_i)$: evidence

Code for H5.10

```
from numpy import log as ln
from math import e
import numpy as np
def fun(x, neg=False):
    if neg:
        return ((e**-x)/(1+e**-x))
    return (1/(1+e**-x))
x = np.array([[1, 2, 2],
             [1, -3, 2],
             [1, -1, 0],
             [1, -2, 2]])
y = np.array([0, 1, 0, 1])
w = np.array([-1, -4, -3])
z = np.array(np.matmul(x, w))
print("z = ", z)
px = np.array([fun(z_i, False) for z_i in z])
print(px)
loss = -1*np.sum((y*ln(1/(1+e**-z)))+(1-y)*ln(e**-z/(1+e**-z)))
print("loss = ", loss)
alpha = 0.6
print("this ", np.sum(x*((y-(1/(1+e**-z)))[::None], axis=0))
nw = w+alpha*np.sum(x*((y-(1/(1+e**-z)))[::None], axis=0)
print("new weights = ", nw)
nz = np.array(np.matmul(x, w))
print("new z", nz)
nloss = -1*np.sum((y*ln(1/(1+e**-nz)))+(1-y)*ln(e**-nz/(1+e**-nz)))
print("new loss", nloss)
```

Code for H5.10

```
https://www.wolframalpha.com/
https://scikit-learn.org/stable/modules/classes.html
https://numpy.org/doc/stable/reference/routines.html
https://pandas.pydata.org/docs/reference/index.html

(x+y)n =  $\sum_{k=0}^n C_k x^n y^k$  =  $\sum_{k=0}^n C_k x^k y^{n-k}$ 

Example: Let us expand  $(x+3)^5$  using the binomial theorem. Here  $y = 3$  and  $n = 5$ . Substituting and expanding, we get:
```

$$(x+3)^5 = 5C_0 x^5 3^0 + 5C_1 x^5-1 3^1 + 5C_2 x^5-2 3^2 + 5C_3 x^5-3 3^3 + 5C_4 x^5-4 3^4 + 5C_5 x^5-5 3^5$$

$$= x^5 + 5x^4 \cdot 3 + 10x^3 \cdot 9 + 10x^2 \cdot 27 + 5x \cdot 81 + 243$$

$$= x^5 + 15x^4 + 90x^3 + 270x^2 + 405x + 243$$

Code for H5.9

```
from numpy import log as ln
from math import e
import numpy as np

def get_stats(confusion_matrix):
    """
    Assumption
    predicted
    1 0
    actual
    1 1 0 TP FN
    0 1 0 FP TN
    """
    TP = confusion_matrix[0, 0]
    TN = confusion_matrix[1, 1]
    FP = confusion_matrix[1, 0]
    FN = confusion_matrix[0, 1]
    # TP = np.diag(confusion_matrix)
    # TN = np.sum(confusion_matrix) - (FP + FN + TP)
```

Sensitivity, hit rate, recall, or true positive rate

TPR = TP/(TP+FN)

Specificity or true negative rate

TNR = TN/(TN+FP)

Precision or positive predictive value

PPV = TP/(TP+FP)

Negative predictive value

NPV = TN/(TN+FN)

Fall out or false positive rate 1 - TNR

FPR = FP/(FP+TN)

False negative rate

FNR = FN/(TP+FN)

False discovery rate

FDR = FP/(TP+FP)

False omission rate 1 - NPV

FOR = FN/(FN+TN)

Overall accuracy

ACC = (TP+TN)/(TP+FP+FN+TN)

actual positivity rate

slightly off in calculations, unsure why

APR = (TP+FN)/(TP+FP+FN+TN)

predicted positivity rate

PPR = (TP+FP)/(TP+FP+FN+TN)

return {'TPR': TPR, 'TNR': TNR, 'PPV': PPV, 'NPV': NPV, 'FPR': FPR, 'FNR': FNR, 'FDR': FDR, 'FOR': FOR, 'ACC': ACC, 'APR': APR, 'PPR': PPR}

a = np.array([13470, 3270], [69, 329])

b = np.array([10049, 3578], [1643, 8692])

c = a+b

T to deal with opposite convention in hw !!!!!

stats_list = [get_stats(a.T), get_stats(b.T), get_stats(c.T)]

for smap in stats_list:

for key, value in smap.items():

print(f'{key} = {value}')

print()

Fairness checks on COMPAS

First, let's compare the accuracy for the two groups:

In other words, we will compute the PPV for each group:

PPV = $\frac{TP}{TP+FP} = P(y=1|\hat{y}=1)$

Again, similar (within a few points). This is a type of fairness known as predictive parity.

We can extend this idea, to check whether a defendant with a given score has the same probability of recidivism for the two groups:

Again, similar (within a few points). This is a type of fairness known as calibration.

We can see that for both African-American and Caucasian defendants, for any given COMPAS score, recidivism rates are similar. This is a type of fairness known as calibration.

Exploratory data analysis - look for possible underlying bias in the data.

Avoid using sensitive group as a feature (or a proxy for a sensitive group), but this doesn't necessarily help. (It didn't help in this example.)

Sometimes, using the sensitive group to explicitly add fairness might be better.

Make sure different groups are well represented in the data.

End users must be made aware of what the model output means -

for example, judges should understand that a "high risk" label only means a 65% chance of reoffending.

Evaluate final model for bias with respect to sensitive groups.

May not be possible to satisfy multiple fairness metrics simultaneously, work with end users to decide which fairness metrics to prioritize, and to create a model that is fair w.r.t. that metric.

xrange=[-1, 1]

```
def generate_polynomial_classifier_data(n=100, xrange=[-1, 1], coefs=[1, 0, 5, 0, 2], sigma=0.5):
    x = np.random.uniform(xrange[0], xrange[1], size=(n, 2))
    ysep = np.polynomial.polynomial.polyval(x[:, 0], coefs)
    y = (x[:, 1]>ysep).astype(int)
```

x[:, 0] = x[:, 0] + sigma * np.random.randn(n)

x[:, 1] = x[:, 1] + sigma * np.random.randn(n)

return x, y

clf = LogisticRegression(penalty='none', tol=0.01, solver='saga')

</div