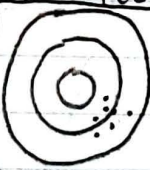Week 5
Model Eval + Refinement

1. • In-sample eval tells how well model will fit the data that's used to train it. However, it can't tell how well the trained model can be used to predict new data.
   • Sol: split data
     - 70% → training
     - 30% → testing (if good, send for training)

   • There's a function in scikit-learn to do this for us.
   1 from sklearn.model_selection import train_test_split
   2 x_train, x_test, y_train, y_test = train_test_split (x_data, y_data, test_size=#, random_state=#)

     - Inputs (vars), ratio of test data, RNG seed, Outputs (arrays)
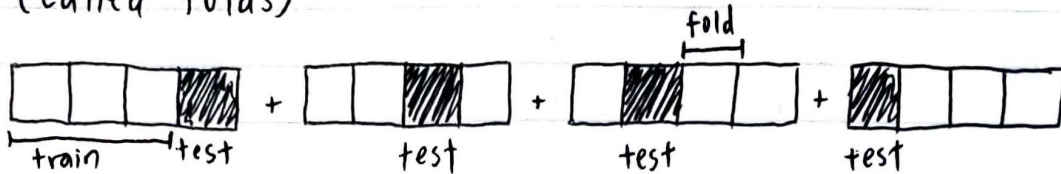
   • Generalization Performance
     - Gen. error is the measure of how well the model does at predicting unseen data
       ▪ Error from testing data is an approximation of this error
     - There is a balance between using training and testing data:

| | ↓ Testing data | ↑ Testing data |
|---|---|---|
| ↓ Training data | Use more data lol | ↓ Accuracy ↑ Precision |
| ↑ Training data | ↑ Accuracy ↓ Precision | Not enough data to go around |

     - To solve this, use cross-validation.

- In cross-validation, data is split into $k$ parts (called folds)


train | test | test | test | test

fold

- Average results are good approximation

- sklearn has a cross-validation function:  — output array of scores

```
1 from sklearn.model_selection import cross_val_score
2 Scores = cross_val_score(lr, x_data, y_data, cv=#)
```
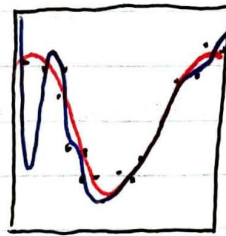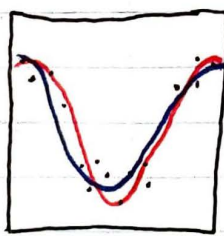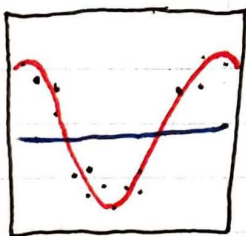Type of model (linreg), # of partitions (folds)
```
3 np.mean(scores)
```

- For more info, cross_val_predict() can be used to generate the predicted y-val of each test element.
```
1 from sklearn.model_selection import cross_val_predict
2 yhat = cross_val_predict(lr2e, x_data, y_data, cv=#)
```

2. How to pick the best polynomial for poly. reg.
- The goal of model selection is to det the order the that's best fit of polynomial y(x).
  - Underfitting: #y < desirable (too simple)
  - Overfitting: #y > desirable (too extreme where few data pts)


Underfitting    Just right    Overfitting (fits noise too well)

Notice overfitting is more "accurate" to training pts but can't extrapolate well


MSE | Test error | Training error | Optimal | Order →

- The order at which $R^2$ is closest to 1 is the optimal order. This can be done in Python:

```
1  Rsqu_test = []   # Init empty list
2  order = [1, 2, 3, 4] # Set to whatever suits you
3  for n in order:
4      pr = PolynomialFeatures(degree = n)
5      x_train_pr = pr.fit_transform(x_train['I.V.'])  # Trans. data
6      x_test_pr = pr.fit_transform(x_test['I.V.'])
7      lr.fit(x_train_pr, y_train)   # Creates model
8      Rsqu_test.append(lr.score(x_test_pr, y_test))  # Calc R²
```

3. Ridge Regression (prev. overfitting + manage outliers)
- Ridge reg. uses a param. "alpha" selected before fitting the data
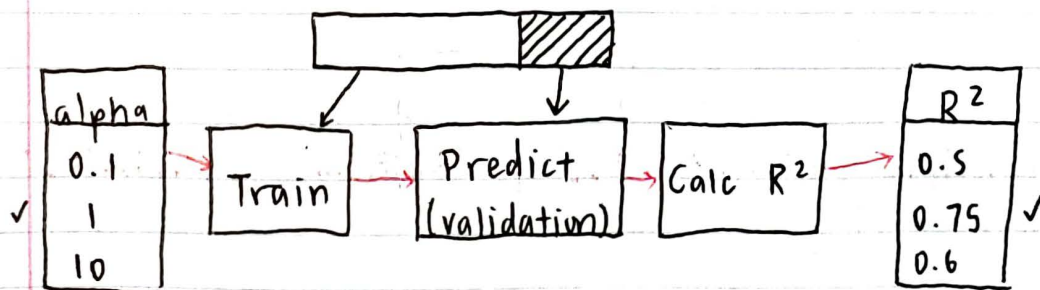- As $\alpha \to \infty$, the coeff.'s $\to 0$ (if $\alpha$ is too high, it can cause underfitting)
- $\alpha$ is selected w/ cross-fitting:

```
1  from sklearn.linear_model import Ridge
2  RidgeModel = Ridge(alpha = #)
3  RidgeModel.fit(x, y)
4  Yhat = RidgeModel.predict(x)
```
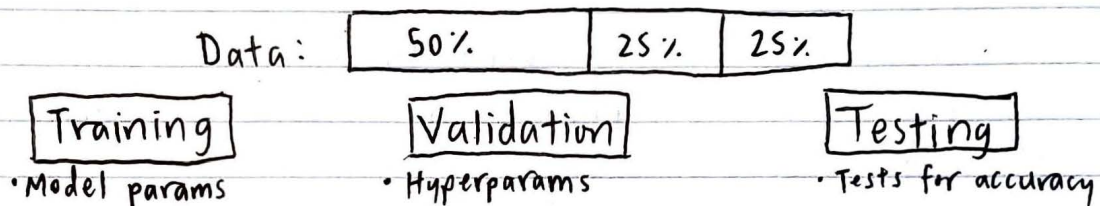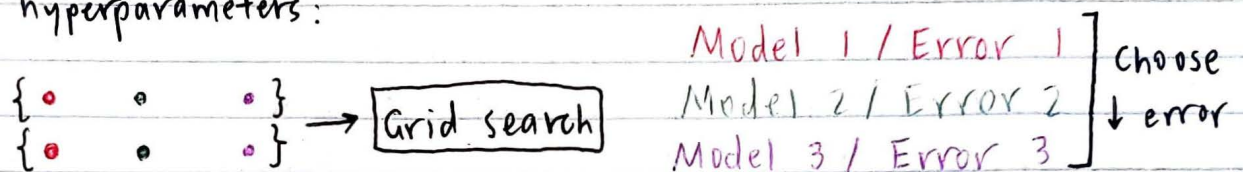
## 4. Grid search

· Terms such as $\alpha$ (not used in training/testing) are called hyperparameters.

 - Scikit-learn has a way to auto-iterate over these hyperparameters w/ cross-validation (known as grid search)

hyperparameters:

$$\begin{Bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{Bmatrix} \rightarrow \boxed{\text{Grid search}}$$

$$\left.\begin{matrix} \text{Model 1 / Error 1} \\ \text{Model 2 / Error 2} \\ \text{Model 3 / Error 3} \end{matrix}\right] \begin{matrix} \text{Choose} \\ \downarrow \text{error} \end{matrix}$$

Data:

| 50% | 25% | 25% |
|-----|-----|-----|
| Training | Validation | Testing |
| ·Model params | · Hyperparams | ·Tests for accuracy |

```
1  from sklearn.linear-model import Ridge
2  from sklearn.model_selection import GridSearchCV
3  parameters1 = [{'alpha': [#1, , #2, ...]}]   *
4  RR = Ridge()
5  Grid1 = GridSearchCV(RR, parameters1, cv=#)
6  Grid1.fit(x_data['var,', 'var2',...],y_data)
7  Grid1.best_estimator_
8  scores = Grid1.cv_results_
9  scores['mean_test_score']
```

For multiple params: ... = [{'alpha': [#1,...], 'normalize': [T, F]}]