# Cheat Sheet – Regularization in ML

## What is Regularization in ML?
- Regularization is an approach to address over-fitting in ML.
- Overfitted model fails to generalize estimations on test data
- When the underlying model to be learned is low bias/high variance, or when we have small amount of data, the estimated model is prone to over-fitting.
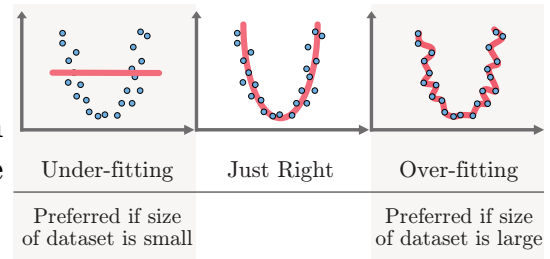- Regularization reduces the variance of the model



Under-fitting | Just Right | Over-fitting

Preferred if size of dataset is small | Preferred if size of dataset is large

Figure 1. Overfitting

## Types of Regularization:

### 1. Modify the loss function:

- **L2 Regularization:** Prevents the weights from getting too large (defined by L2 norm). Larger the weights, more complex the model is, more chances of overfitting.

$$loss = error(y, \hat{y}) \boxed{+ \lambda \sum_j \beta_j^2} \quad \lambda \geq 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

- **L1 Regularization:** Prevents the weights from getting too large (defined by L1 norm). Larger the weights, more complex the model is, more chances of overfitting. L1 regularization introduces sparsity in the weights. It forces more weights to be zero, than reducing the the average magnitude of all weights

$$loss = error(y, \hat{y}) \boxed{+ \lambda \sum_j |\beta_j|} \quad \lambda \geq 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

- **Entropy:** Used for the models that output probability. Forces the probability distribution towards uniform distribution.

$$loss = error(p, \hat{p}) \boxed{- \lambda \sum_i \hat{p}_i log(\hat{p}_i)} \quad \lambda \geq 0, \ \lambda \propto model \ bias, \ \lambda \propto \frac{1}{model \ variance}$$

### 2. Modify data sampling:
- **Data augmentation:** Create more data from available data by randomly cropping, dilating, rotating, adding small amount of noise etc.
- **K-fold Cross-validation:** Divide the data into k groups. Train on (k-1) groups and test on 1 group. Try all k possible combinations.

### 3. Change training approach:
- **Injecting noise:** Add random noise to the weights when they are being learned. It pushes the model to be relatively insensitive to small variations in the weights, hence regularization
- **Dropout:** Generally used for neural networks. Connections between consecutive layers are randomly dropped based on a dropout-ratio and the remaining network is trained in the current iteration. In the next iteration, another set of random connections are dropped.
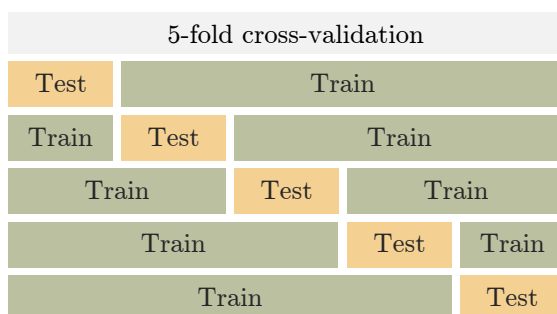


5-fold cross-validation
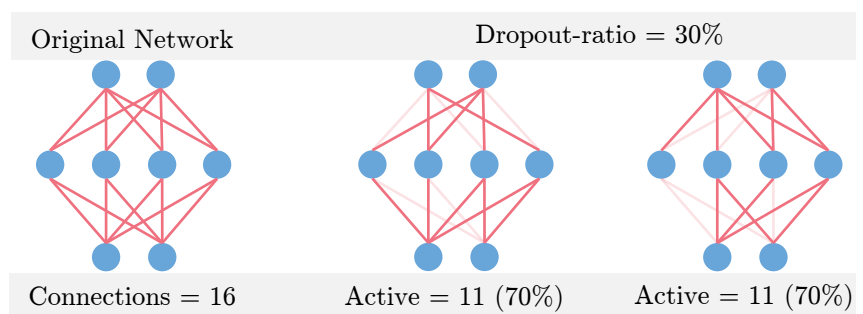
| Test | Train | | |
| Train | Test | Train | |
| Train | Test | Train | |
| Train | Test | Train |
| Train | Test |

Figure 2. K-fold CV



Original Network | Dropout-ratio = 30%

Connections = 16 | Active = 11 (70%) | Active = 11 (70%)

Figure 3. Drop-out