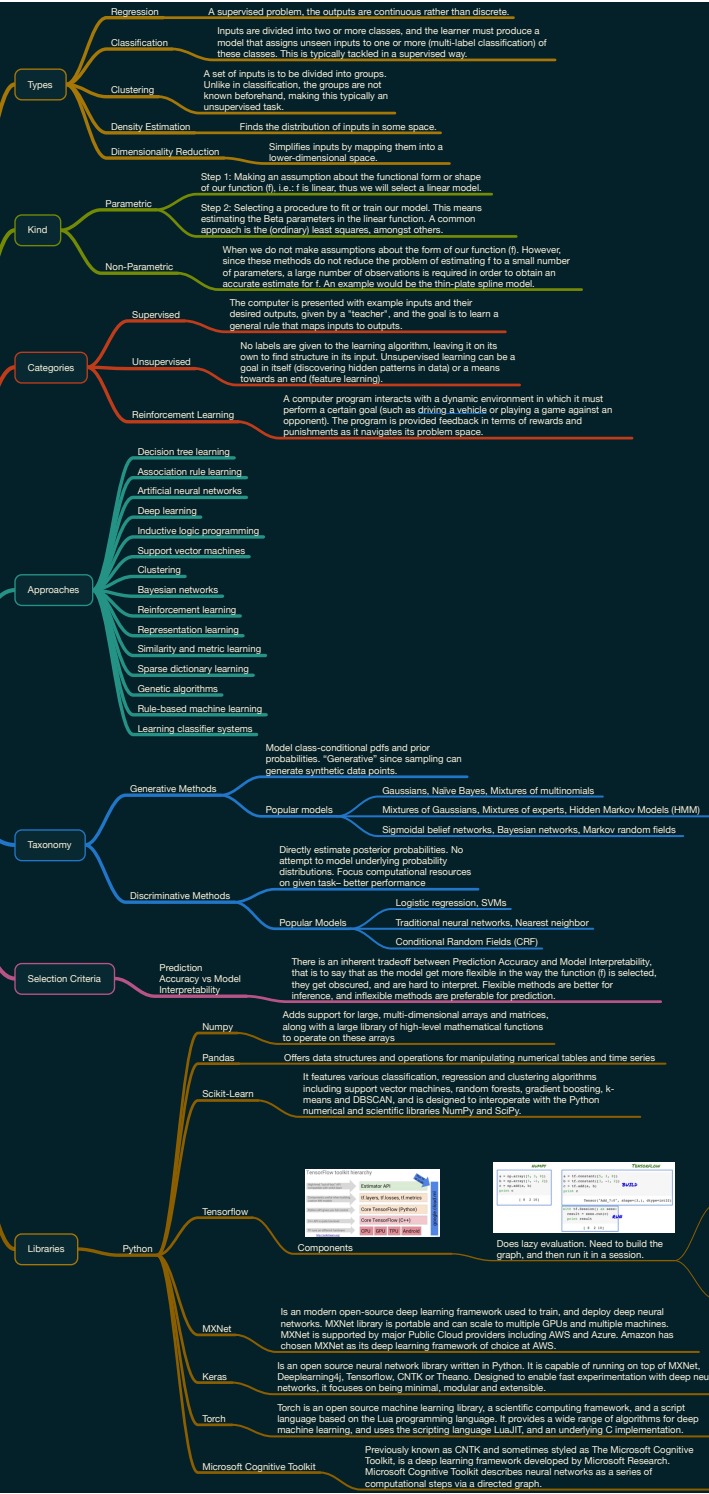


# Machine Learning Concepts



**Motivation**

- Prediction**: When we are interested mainly in the predicted variable as a result of the inputs, but not on the each way of the inputs affect the prediction. In a real estate example, Prediction would answer the question of: Is my house over or under valued? Non-linear models are very good at these sort of predictions, but not great for inference because the models are much less interpretable.
- Inference**: When we are interested in the way each one of the inputs affect the prediction. In a real estate example, Inference would answer the question of: How much would my house cost if I had a view of the sea? Linear models are more suited for inference because the models themselves are easier to understand than their non-linear counterparts.

	Actual Positive	Actual Negative	Total
Predicted Positive	90 (90%)	10 (10%)	100
Predicted Negative	10 (10%)	90 (90%)	100
Total	100	100	200

**Confusion Matrix**

Fraction of correct predictions, not reliable as skewed when the data set is unbalanced (that is, when the number of samples in different classes vary greatly)

**Precision**

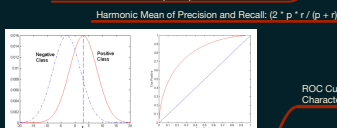
$(TP / (TP + FP))$   
Which tells us what proportion of patients we diagnosed as having cancer actually had cancer. In other words, proportion of TP in the set of positive cancer diagnoses. This is given by the rightmost column in the confusion matrix.

**Recall**

$(TP / (TP + FN))$   
Which tells us what proportion of patients that actually had cancer were diagnosed by us as having cancer. In other words, proportion of TP in the set of real cancer cases. This is given by the bottom row in the confusion matrix.

**F1 score**

Out of all the examples the classifier labeled as positive, what fraction were correct?  
Out of all the positive examples there were, what fraction did the classifier pick up?  
Harmonic Mean of Precision and Recall:  $2 * p * r / (p + r)$



**Bias-Variance Tradeoff**

Bias refers to the amount of error that is introduced by approximating a real-life problem, which may be extremely complicated, by a simple model. If Bias is high, and/or if the algorithm performs poorly even on your training data, try adding more features, or a more flexible model.

**Goodness of Fit = R^2**

Variance is the amount our model's prediction would change when using a different training data set. High: Remove features, or obtain more data.

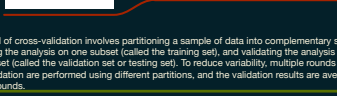
**Mean Squared Error (MSE)**

$1.0 - \text{sum of squared errors} / \text{total sum of squares}$   
 $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$

**Error Rate**

The proportion of mistakes made if we apply our estimate model function to the training observations in a classification setting.

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$



**Grid Search**

The traditional way of performing hyperparameter optimization has been grid search, or a parameter sweep, which is simply an exhaustive searching through a manually specified subset of the hyperparameter space of a learning algorithm. A grid search algorithm must be guided by some performance metric, typically measured by cross-validation on the training set or evaluation on a held-out validation set.

**Random Search**

Since grid searching is an exhaustive and therefore potentially expensive method, several alternatives have been proposed. In particular, a randomized search that simply samples parameter settings a fixed number of times has been found to be more effective in high-dimensional spaces than exhaustive search.

**Gradient-based optimization**

For specific learning algorithms, it is possible to compute the gradient with respect to hyperparameters and then optimize the hyperparameters using gradient descent. The first usage of these techniques was focused on neural networks. Since then, these methods have been extended to other models such as support vector machines or logistic regression.

**Early Stopping (Regularization)**

Early stopping rules provide guidance as to how many iterations can be run before the learner begins to over-fit, and stop the algorithm then.

**Overfitting**

When a given method yields a small training MSE (or cost), but a large test MSE (or cost), we are said to be overfitting the data. This happens because our statistical learning procedure is trying too hard to find patterns in the data, that might be due to random chance, rather than a property of our function. In other words, the algorithm may be learning the training data too well. If model overfits, try removing some features, decreasing degrees of freedom, or adding more data.

**Underfitting**

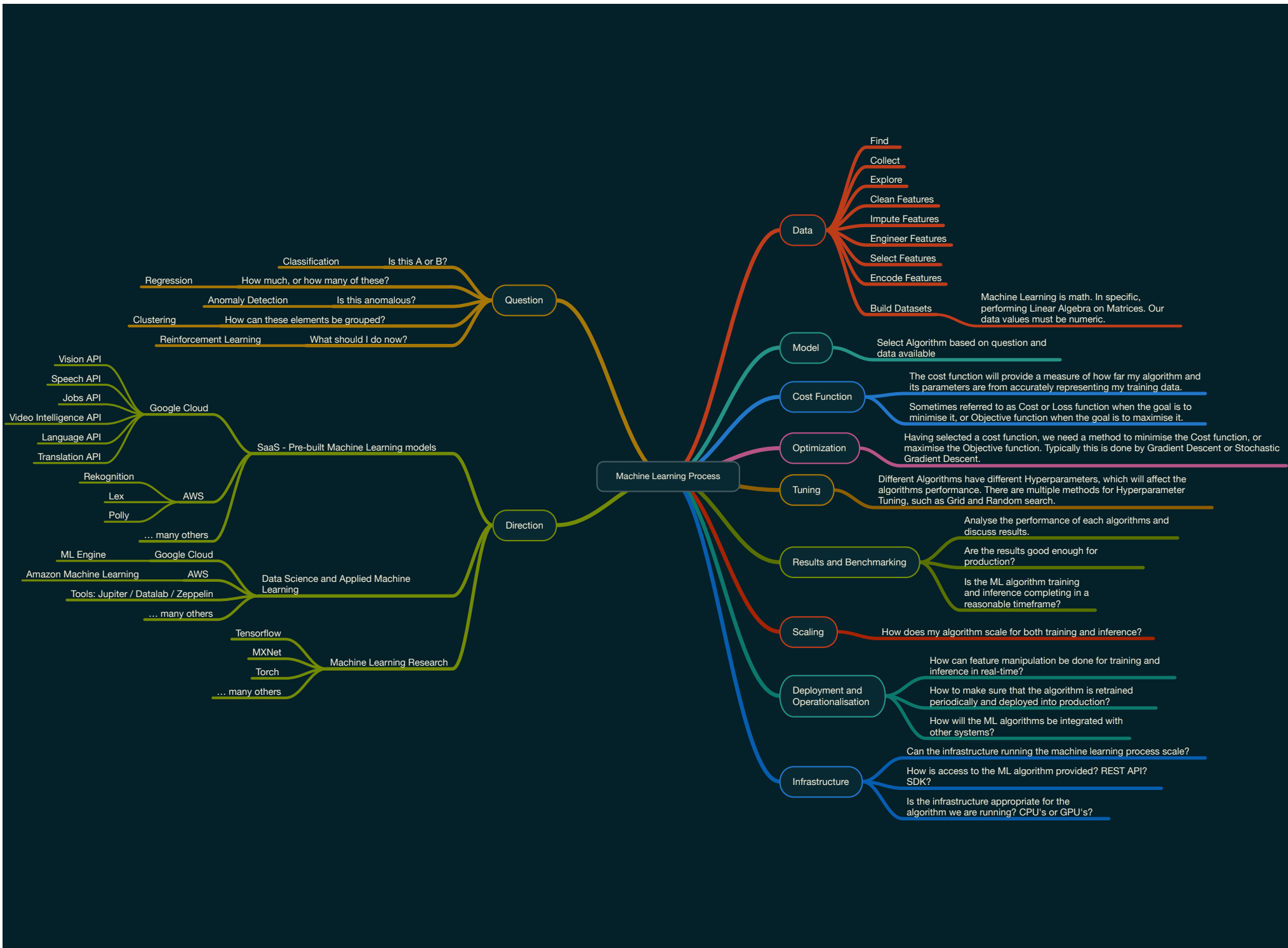
Opposite of Overfitting. Underfitting occurs when a statistical model or machine learning algorithm cannot capture the underlying trend of the data. It occurs when the model or algorithm does not fit the data enough.

**Bootstrap**

Underfitting occurs if the model or algorithm shows low variance but high bias (to contrast the opposite, overfitting from high variance and low bias). It is often a result of an excessively simple model. Test that applies Random Sampling with Replacement of the available data, and assigns measures of accuracy (bias, variance, etc.) to sample estimates.

**Bagging**

An approach to ensemble learning that is based on bootstrapping. Shortly, given a training set, we produce multiple different training sets (called bootstrap samples), by sampling with replacement from the original dataset. Then, for each bootstrap sample, we build a model. The results in an ensemble of models, where each model votes with the equal weight. Typically, the goal of this procedure is to reduce the variance of the model of interest (e.g. decision trees).

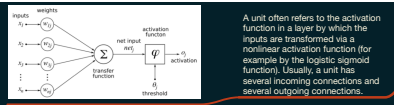




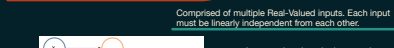
# Machine Learning Models

## Neural Networks

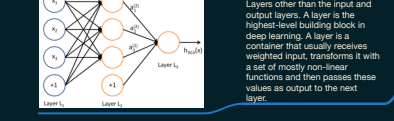
### Unit (Neurons)



### Input Layer



### Hidden Layers



### Batch Normalization

With SGD, the training proceeds in steps, and at each step we consider a mini-batch  $x_i$  of size  $m$ . The mini-batch is used to approximate the gradient of the loss function with respect to the parameters.

Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch is an estimate of the gradient over the training set, whose quality improves as the batch size increases. Second, computation over a batch can be much more efficient than in computations for individual examples, due to the parallelism afforded by the modern computing platforms.

### Learning Rate

However, if you actually try that, the weights will change far too much each iteration, which will make them "overcorrect" and the loss will actually increase/diverge. So in practice, people usually multiply each derivative by a small value called the "learning rate" before they subtract it from its corresponding weight.

Neural networks are often trained by gradient descent on the weights. This means at each iteration we use backpropagation to calculate the derivative of the loss function with respect to each weight and subtract it from that weight.

Simplest recipe: keep it fixed and use the same for all parameters.

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_L(\theta)$$

### Tricks

Reduce by 0.5 when validation error stops improving

$$\alpha = \frac{\alpha_0 \tau}{\max(t, \tau)}$$

Reduction by  $O(1/\tau)$  because of theoretical convergence guarantees, with hyper-parameters  $\epsilon_0$  and  $\tau$  and  $t$  is iteration numbers

Better results by allowing learning rates to decrease Options:

But, this turns out to be a mistake, because if every neuron in the network computes the same output, then they will also all compute the same gradients during back-propagation and undergo the exact same parameter updates. In other words, there is no source of asymmetry between neurons if their weights are initialized to be the same.

In the ideal situation, with proper data normalization it is reasonable to assume that approximately half of the weights will be positive and half of them will be negative. A reasonable-sounding idea then might be to set all the initial weights to zero, which you expect to be the "best guess" in expectation.

### Weight Initialization

All Zero Initialization

Initialization with Small Random Numbers

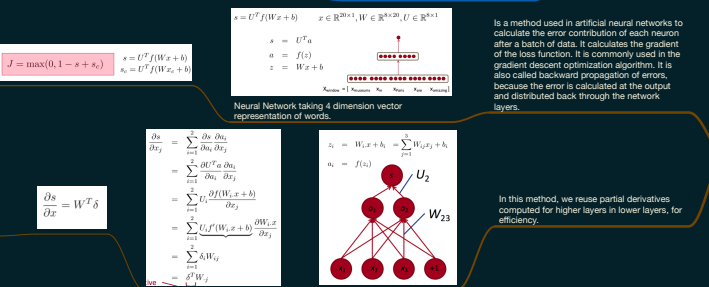
Thus, you still want the weights to be very close to zero, but not identically zero. In this way, you can random these neurons to small numbers which are very close to zero, and it is also possible to use small numbers drawn from a uniform distribution, but this seems to have relatively little impact on the final performance in practice.

One problem with the above suggestion is that the distribution of the outputs from a randomly initialized neuron has a variance that grows with the number of inputs. It turns out that you can normalize the variance of each neuron's output to 1, by scaling its weight vector by the square root of its fan-in (i.e., its number of inputs)

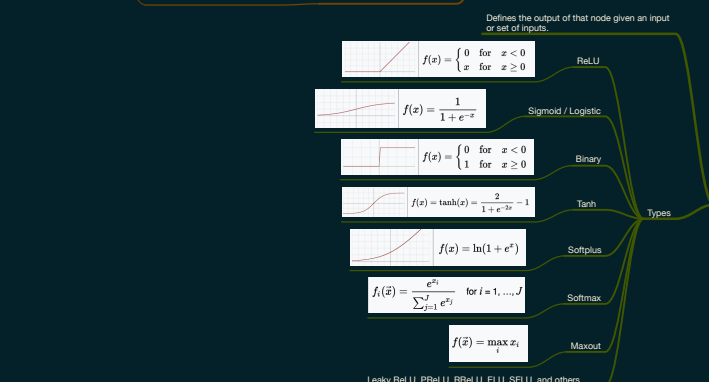
Calibrating the Variances

This ensures that all neurons in the network initially have approximately the same output distribution and empirically improves the rate of convergence. The detailed derivations can be found from Page. 18 to 23 of the slides. Please note that, in the derivations, it does not consider the influence of ReLU neurons.

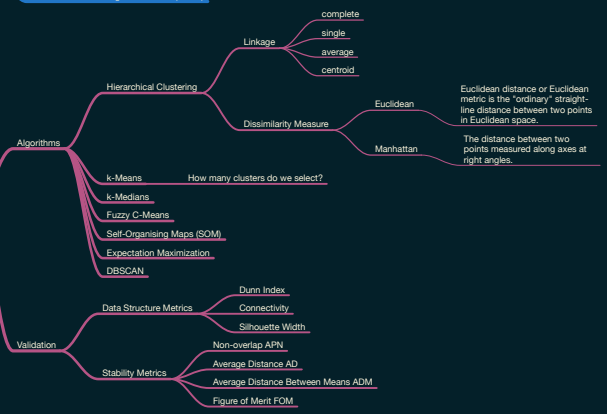
### Backpropagation



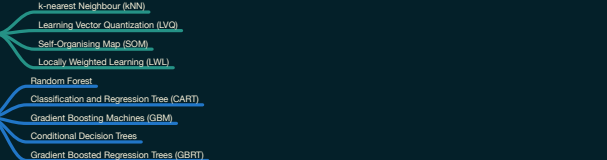
### Activation Functions



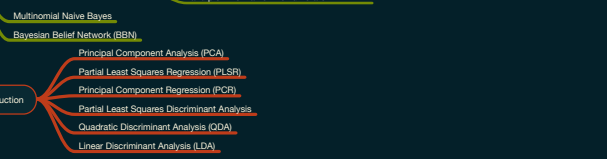
### Clustering



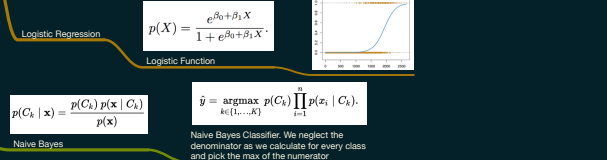
### Decision Tree



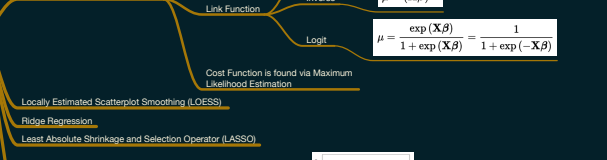
### Instance Based



### Dimensionality Reduction



### Bayesian



### Regression

