

## Assignment 2: Guest Lecturer Edition

---

Aditya Shinde

September 14, 2017

### Questions

#### LINEAR REGRESSION [20PTS]

For this question, I watched a couple of videos on youtube about maximum likelihood estimation and sampling data points from probability distributions.

1

If the error terms  $\epsilon_i$  are assumed to be i.i.d and sampled from a Gaussian distribution,

$$\begin{aligned}\epsilon &\sim N(0, \sigma^2) \\ \epsilon &= y - w^T x \\ y &\sim N(w^T x, \sigma^2)\end{aligned}$$

Similarly, If

$$\begin{aligned}\epsilon &\sim Lap(0, b) \\ \epsilon &= y - w^T x \\ y &\sim Lap(w^T x, b)\end{aligned}$$

$$P(Y|x_i; w) = \operatorname{argmax}_i \Pi_i \operatorname{Lap}(w^T x, b)$$

$$P(Y|x_i; w) = \operatorname{argmax}_i \frac{1}{2b} e^{\left(\frac{-|y - w^T x|}{b}\right)}$$

$$\log P(Y|x_i; w) = \operatorname{argmax}_i \sum -\log 2b + \left(\frac{-|y - w^T x|}{b}\right)$$

Now,  $\log 2b$  is a constant,

So we can write  $J_{Lap}(w)$  as,

$$J_{Lap}(w) = \operatorname{argmin}_w \sum (|y - w^T x|)$$

2

If the noise terms are distributed as Gaussians, the loss function is in the form of squared errors. This places a very high penalty on outliers and thus the model is highly reactive towards them. In case of a Laplacian distribution, the loss function is the mean absolute error. This is less reactive to outliers.

## REGULARIZATION [30PTS]

1

As  $\lambda$  decreases, the cost function penalises exploding weights to a lesser extent.

As  $\lambda \rightarrow 0$  the cost function  $J_R(\beta) \rightarrow J(\beta)$

As  $\lambda$  increases, the regularization term becomes more dominant. It starts heavily penalising all weights. As  $\lambda \rightarrow \infty$ , the cost function  $J_R(\beta) \sim \lambda ||\beta||^2$

The model will only fit the intercept.

2

$$\beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \prod_{i=1}^n P(Y_i|X_i; \beta) P(\beta)$$

$$\log \beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \sum_{i=1}^n (\log P(Y_i|X_i; \beta) + \log P(\beta))$$

$$\log \beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \sum_{i=1}^n \left( -(Y - W^T X)^2 + \left( -\log \frac{\lambda}{I\sigma^2} - \frac{\beta^2 \lambda}{I\sigma^2} \right) \right)$$

Getting rid of the constants.

$$\operatorname{argmax}_{\beta} \log \beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \sum_{i=1}^n - \left( (Y - W^T X)^2 + \frac{\beta^2 \lambda}{I\sigma^2} \right)$$

If we assume unit variance,  $\sigma^2 = 1$  and  $I$  is an identity matrix. And multiply and divide right hand side by -1 to change argmax to argmin.

$$\operatorname{argmax}_{\beta} \log \beta_{\text{MAP}} = \operatorname{argmin}_{\beta} \sum_{i=1}^n ((Y - W^T X)^2 + \lambda \beta^2)$$

3

Now in this case  $\beta$  is a random variable with variance  $\frac{I\sigma^2}{\lambda}$ .

If  $\lambda \rightarrow 0$  in this case, that means  $\operatorname{Var}[\beta] \rightarrow \infty$ , that is variance tends to infinity. So the samples will be far away.

Similarly, if  $\lambda \rightarrow \infty$ ,  $\operatorname{Var}[\beta] \rightarrow 0$  which means  $p(\beta)$  will be a single point at 0 and we would essentially be fitting only the intercept.

## EVOLUTIONARY COMPUTATION

1

Generational genetic algorithms use specific distinct generations as different sets of parameters. By keeping distinct generations across epochs, we make sure that none of the individuals from the previous generation are in the current generation. This is specially helpful when the cost function is dynamic. It can also be helpful to find multiple solutions.

Steady state GAs on the other hand use a replacement mechanism. The newer individuals which are generated using genetic operators are swapped with the weaker individuals of the current generation.

One of the advantages of steady state GA's is that they require less resources in terms of memory since the population size is constant. Generational GA's on the other hand need twice as much memory. But steady state GA's usually take more time to converge.

2

Pittsburgh and Michigan classifiers both work on the same principles of evolutionary computation. The distinction is in the role of a single individual. In the Pittsburgh classifiers, each individual is a complete solution consisting of many rules. In case Michigan classifiers, an individual is just an element of a rule based system. Many such individuals need to be combined to make a complete solution.

One of the advantages that Michigan classifiers have is diversity. Since the individuals are made of elementary rules, they can be used with many problems easily as compared to Pittsburgh.