

Midterm Practice Problems Solutions

Problem 1 (DECISION TREE)

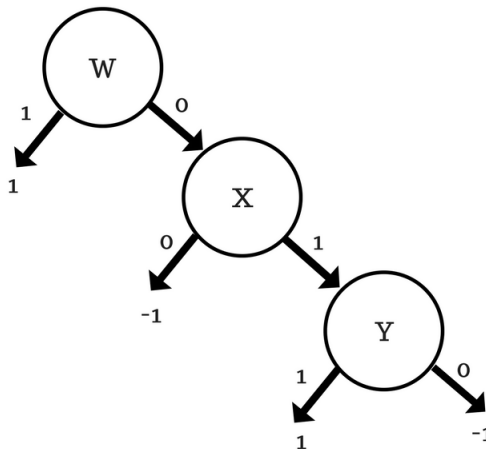
(a) $H(Z) = H(\frac{2}{5}) = \frac{2}{5} \log(\frac{5}{2}) + \frac{3}{5} \log(\frac{5}{3}) = 0.970951$

$$\begin{aligned} H(Z|split(W)) &= \frac{1}{5}H(Z|branch(W = 1)) + \frac{4}{5}H(Z|branch(W = 0)) \\ &= \frac{1}{5}H(1) + \frac{4}{5}H(\frac{1}{4}) = 0.649 \end{aligned}$$

$$\begin{aligned} H(Z|split(X)) &= \frac{2}{5}H(Z|branch(X = 1)) + \frac{3}{5}H(Z|branch(X = 0)) \\ &= \frac{2}{5}H(0.5) + \frac{3}{5}H(\frac{1}{3}) = 0.951 \end{aligned}$$

$$\begin{aligned} H(Z|split(Y)) &= \frac{2}{5}H(Z|branch(Y = 1)) + \frac{3}{5}H(Z|branch(Y = 0)) \\ &= \frac{2}{5}H(0.5) + \frac{3}{5}H(\frac{1}{3}) = 0.951 \end{aligned}$$

Therefore splitting with W has the highest information gain of 0.321951 in comparison to splitting with X or Y both having IG of 0.0199



(b)

X and Y can be swapped to get another possible decision tree.

- (c) Samples# 5 and 7 have the same features but different labels, therefore it is not possible to construct such a decision tree.

Problem 2 (PERCEPTRON)

- (a) AND

$\theta = (2, 2, -3)$ if the augmented features are $(x_1, x_2, 1)$. Multiple solutions are possible.

- (b) XOR

No solution exists because the data is not linearly separable.

Problem 3 (LOGISTIC REGRESSION)

- (a) For finding $\nabla_{\mathbf{w}} J(\mathbf{w}, \mathbf{x})$ first note that $\alpha_i(\mathbf{x})$ doesn't depend on \mathbf{w} . Therefore the derivation will be similar to that for logistic regression except for the additional localized weights.

$$\nabla_{\mathbf{w}} J(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^n \alpha_i(\mathbf{x})(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)\mathbf{x}_i.$$

Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$, then we can write $\nabla_{\mathbf{w}} J(\mathbf{w}, \mathbf{x}) = \mathbf{X}^T \mathbf{z}$, where \mathbf{z} is a vector of \mathbb{R}^n with each entree $\alpha_i(\mathbf{x})(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)$.

- (b) From the expression for the gradient you can see that

$$\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial w_j} = \sum_{i=1}^n \alpha_i(\mathbf{x})(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)x_{i,j}$$

$$\begin{aligned} \frac{\partial^2 J(\mathbf{w}, \mathbf{x})}{\partial w_k \partial w_j} &= \frac{\partial}{\partial w_k} \left(\sum_{i=1}^n \alpha_i(\mathbf{x})(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)x_{i,j} \right) \\ &= \sum_{i=1}^n \alpha_i(\mathbf{x}) \frac{\partial}{\partial w_k} h_{\mathbf{w}}(\mathbf{x}_i)x_{i,j} \\ &= \sum_{i=1}^n \alpha_i(\mathbf{x}) h_{\mathbf{w}}(\mathbf{x}_i)(1 - h_{\mathbf{w}}(\mathbf{x}_i))x_{i,j}x_{i,k} \end{aligned}$$

Therefore we have

$$\nabla_{\mathbf{w}}^2 J(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^n \alpha_i(\mathbf{x}) h_{\mathbf{w}}(\mathbf{x}_i)(1 - h_{\mathbf{w}}(\mathbf{x}_i))\mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}^T \mathbf{D} \mathbf{X}$$

$$\mathbf{u}^T \mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{u} = \|D^{\frac{1}{2}} \mathbf{X} \mathbf{u}\|_2^2 > 0 \quad \forall \mathbf{u} \neq 0$$

- (c)

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta \mathbf{X}^T \mathbf{z}^t$$

- (d) Non parametric method.

Problem 4 (LINEAR REGRESSION)

(a)

(b) Gradient Descent:

$$w_1^{t+1} \leftarrow w_1^t + \eta \sum_{i=1}^M (y_i - (w_1^t x_1^{(i)} + w_2^t x_2^{(i)})) x_1^{(i)} \quad (1)$$

$$w_2^{t+1} \leftarrow w_2^t + \eta \sum_{i=1}^M (y_i - (w_1^t x_1^{(i)} + w_2^t x_2^{(i)})) x_2^{(i)} \quad (2)$$

You can also do stochastic gradient descent.

(c) To prove Eq. (??) has a global optimum, we have to show the function is convex. To prove the function is convex, we need to demonstrate the Hessian matrix (or the second derivatives) is positive semi-definite.

The Hessian of Eq. (??) is

$$H = \begin{bmatrix} \sum (x_1^{(i)})^2 & \sum x_1^{(i)} x_2^{(i)} \\ \sum x_1^{(i)} x_2^{(i)} & \sum (x_2^{(i)})^2 \end{bmatrix} \quad (3)$$

To show H is positive semi-definite, we have to prove for every vector $z \neq 0$, $z^T H z \geq 0$. This can be done by the following equations:

$$\begin{aligned} z^T H z &= \sum (x_1^{(i)})^2 z_1^2 + 2 \sum x_1^{(i)} x_2^{(i)} z_1 z_2 + \sum (x_2^{(i)})^2 z_2^2 \\ &= \sum (z_1 x_1^{(i)} + z_2 x_2^{(i)})^2 \\ &\geq 0 \end{aligned} \quad (4)$$

Problem 5 (MAXIMUM LIKELIHOOD ESTIMATION)

(a)

$$\begin{aligned} \mathcal{L}(x_1, x_2, \dots, x_n; \lambda) &= \prod_{i=1}^n f(x_i | \lambda) \\ &= \lambda^n (x_1 x_2 \dots x_n)^{-\lambda-1} \end{aligned}$$

(b)

$$\begin{aligned} \arg \max_{\lambda} \mathcal{L}(x_1, x_2, \dots, x_n; \lambda) &= \arg \min_{\lambda} \left(-\ln \mathcal{L}(x_1, x_2, \dots, x_n; \lambda) \right) \\ -\mathcal{L}\mathcal{L} &= -n \ln \lambda + (\lambda + 1) \sum_{i=1}^n \ln x_i \\ \frac{d(-\mathcal{L}\mathcal{L})}{d\lambda} &= -\frac{n}{\lambda} + \sum_{i=1}^n \ln x_i = 0 \\ \lambda_{mle} &= \frac{n}{\sum_{i=1}^n \ln x_i} \end{aligned}$$

But since $\lambda > 1$ is known about the model, $\lambda_{mle} = \max(1, \frac{n}{\sum_{i=1}^n \ln x_i})$

(c) $\lambda_{mle} = \max(1, \frac{4}{6}) = 1$

Problem 6 (MAXIMUM LIKELIHOOD ESTIMATION 2)

(a)

$$\begin{aligned} l(\theta) &= \sum_n \log P(X_i; \theta) \\ &= \sum_n x_n \log(\theta) + (1 - x_n) \log(1 - \theta) \\ &= 3 \log(\theta) + \log(1 - \theta) \end{aligned}$$

(b)

$$l'(\theta) = \frac{3}{\theta} - \frac{1}{1 - \theta}$$

(c)

$$\begin{aligned} l'(\theta) &= \frac{3}{\theta} - \frac{1}{1 - \theta} = 0 \\ \hat{\theta} &= \frac{3}{4} \end{aligned}$$

Problem 7 (KERNEL)

To show that K_{β} is a kernel, simply use the same feature mapping as used by the polynomial kernel of degree 3, but first scale \mathbf{x} by $\sqrt{\beta}$. So, in effect they are both polynomial kernels of degree 3. If you look at the resulting feature vector, the offset term 1 is unchanged, the linear terms are scaled by $\sqrt{\beta}$, the quadratic terms are scaled by β , and the cubic terms are scaled by $\beta^{1.5}$. Although the

model class remains unchanged, this changes how we penalize the features during learning (from the $\|\theta\|^2$ in the objective). In particular, higher-order features will become more costly to use, so this will bias more towards a lower-order polynomial.

That is,

$$\begin{aligned}
 K_\beta(\mathbf{x}, \mathbf{z}) &= (1 + \beta \mathbf{x} \cdot \mathbf{z})^3 \\
 &= (1 + \beta (x_1 z_1 + x_2 z_2))^3 \\
 &= 1 + 3\beta (x_1 z_1 + x_2 z_2) + 3\beta^2 (x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2) \\
 &\quad + \beta^3 (x_1^3 z_1^3 + 3x_1^2 z_1^2 x_2 z_2 + 3x_1 z_1 x_2^2 z_2^2 + x_2^3 z_2^3)
 \end{aligned}$$

so that

$$\phi_\beta(\mathbf{x}) = (1, \sqrt{3\beta}x_1, \sqrt{3\beta}x_2, \sqrt{3\beta}x_1^2, \sqrt{6\beta}x_1x_2, \sqrt{3\beta}x_2^2, \sqrt{\beta^3}x_1^3, \sqrt{3\beta^3}x_1^2x_2, \sqrt{3\beta^3}x_1x_2^2, \sqrt{\beta^3}x_2^3)^T$$

The m^{th} -order terms in $\phi_\beta(\cdot)$ are scaled by $\beta^{m/2}$, so β trades off the influence of the higher-order versus lower-order terms in the polynomial. If $\beta = 1$, then $\beta^{1/2} = \beta = \beta^{3/2}$ so that $K_\beta = K$. If $0 < \beta < 1$, then $\beta^{1/2} > \beta > \beta^{3/2}$ so that lower-order terms have more weight and higher-order terms less weight; as $\beta \rightarrow 0$, K_β approaches $1 + 3\beta \mathbf{x} \cdot \mathbf{z}$ (a linear separator). If $\beta > 1$, the trade-off is reversed; as $\beta \rightarrow \infty$, only the constant and cubic terms in K_β remain.