

Student Name: Deeksha Arora

Roll Number: 20111017

Date: October 30, 2020

The given absolute loss regression problem with l_1 regularization is

$$\mathbf{w}_{opt} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n| + \lambda \|\mathbf{w}\|_1, \text{ where } \|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$$

Let $L_1(\mathbf{w}) = \sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n|$ and $L_2(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$

The function $L_1(\mathbf{w})$ is summation of N absolute value functions. We know that absolute value function is always convex and second derivative of convex functions is positive. So, the second derivative of sum of convex functions is also positive. Since, second derivative of sum of convex functions is positive, therefore, sum of two convex functions is also a convex function. Hence, function $L_1(\mathbf{w})$ is convex in nature.

The function $L_2(\mathbf{w})$ is absolute value function and hence convex.

Let $L(\mathbf{w})$ represents the regularized loss function, therefore

$$L(\mathbf{w}) = \sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n| + \lambda \|\mathbf{w}\|_1$$

$$\Rightarrow L(\mathbf{w}) = L_1(\mathbf{w}) + L_2(\mathbf{w})$$

Since $L_1(\mathbf{w})$ and $L_2(\mathbf{w})$ are convex functions and we know that sum of two convex functions is also convex, so $L(\mathbf{w})$ is a **convex function**.

Even though absolute value function is convex, it is non-differentiable. So we define subgradient for $L(\mathbf{w})$ as:

$$\partial_{\mathbf{w}} L(\mathbf{w}) = \partial_{\mathbf{w}} \left(\sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n| + \lambda \|\mathbf{w}\|_1 \right)$$

From Sum Rule of Differentiation:

$$\partial_{\mathbf{w}} L(\mathbf{w}) = \partial_{\mathbf{w}} L_1(\mathbf{w}) + \partial_{\mathbf{w}} L_2(\mathbf{w})$$

Sub Gradient of $L_1(\mathbf{w}) = \sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n|$:

Let $l_n(\mathbf{w}) = |y_n - \mathbf{w}^\top \mathbf{x}_n|$. This function is non-differentiable at $|y_n - \mathbf{w}^\top \mathbf{x}_n| = 0$.

Let \mathbf{u}'_n be the sub-gradient vector for $l_n(\mathbf{w})$, defined as:

$$\mathbf{u}'_n = \partial_w (l_n(\mathbf{w}))$$

$$\mathbf{u}'_n = \begin{cases} \mathbf{x}_n & \text{if } y_n - \mathbf{w}^\top \mathbf{x}_n < 0, \\ -c\mathbf{x}_n, \text{ where } c \in [-1, +1] & \text{if } y_n - \mathbf{w}^\top \mathbf{x}_n = 0 \\ -\mathbf{x}_n & \text{if } y_n - \mathbf{w}^\top \mathbf{x}_n > 0, \end{cases} \quad (1)$$

Let \mathbf{u} is the sub-gradient vector for $L_1(\mathbf{w}) = \sum_{n=1}^N |y_n - \mathbf{w}^\top \mathbf{x}_n|$, defined as:

$$\begin{aligned} \mathbf{u} &= \partial_w L_1(\mathbf{w}) \\ &= \sum_{n=1}^N \partial_w (l_n(\mathbf{w})) \\ &= \sum_{n=1}^N \mathbf{u}'_n \end{aligned} \quad (2)$$

where values of \mathbf{u}'_n can be substituted from equation (1).

Sub-Gradient of $L_2(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$, where $\|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$:

It is absolute value function and is non-differentiable at 0 and hence we define subgradient for it. Let \mathbf{v} be the sub-gradient vector for $L_2(\mathbf{w})$, then the i^{th} element of this sub gradient vector is given by:

$$v_i = \partial_{w_i} (L_2(\mathbf{w})) = \begin{cases} -\lambda, & \text{if } w_i < 0, \\ [-\lambda, +\lambda], & \text{if } w_i = 0 \\ \lambda, & \text{if } w_i > 0, \end{cases} \quad (3)$$

Sub-Gradient of $L(\mathbf{w}) = L_1(\mathbf{w}) + L_2(\mathbf{w})$:

Let \mathbf{g} denotes the sub-gradient vector of $L(\mathbf{w})$, then the i^{th} element of this vector is given by:

$$g_i = u_i + v_i \quad (4)$$

The values of u_i and v_i are substituted from equations (2) and (3) respectively.

Student Name: Deeksha Arora

Roll Number: 20111017

Date: October 30, 2020

The new squared loss function with masked inputs is defined as :

$$\begin{aligned} L(\mathbf{w}) &= \sum_{n=1}^N (y_n - \mathbf{w}^T \tilde{\mathbf{x}}_n)^2, \text{ where } \tilde{\mathbf{x}}_n = \mathbf{x}_n \circ \mathbf{m}_n \\ &= \sum_{n=1}^N (y_n^2 - y_n \mathbf{w}^T \tilde{\mathbf{x}}_n - y_n \mathbf{w} \tilde{\mathbf{x}}_n^T + \mathbf{w} \tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n \mathbf{w}^T) \end{aligned} \quad (5)$$

Taking the expectation w.r.t. \mathbf{m}_n we get:

$$\begin{aligned} E[L(\mathbf{w})] &= \sum_{n=1}^N E[y_n^2 - y_n \mathbf{w}^T \tilde{\mathbf{x}}_n - y_n \mathbf{w} \tilde{\mathbf{x}}_n^T + \mathbf{w} \tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n \mathbf{w}^T] \\ &= \sum_{n=1}^N E[y_n^2] - E[y_n \mathbf{w}^T \tilde{\mathbf{x}}_n] - E[y_n \mathbf{w} \tilde{\mathbf{x}}_n^T] + E[\mathbf{w} \tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n \mathbf{w}^T] \\ &= \sum_{n=1}^N y_n^2 - y_n \mathbf{w}^T E[\tilde{\mathbf{x}}_n] - y_n \mathbf{w} E[\tilde{\mathbf{x}}_n^T] + \mathbf{w} E[\tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n] \mathbf{w}^T \end{aligned} \quad (6)$$

The expected value of a matrix is the matrix of cell-wise expected values, so

$$E[\tilde{\mathbf{x}}]_i = E[(\mathbf{x} \circ \mathbf{m})_i] = \mathbf{x}_i E[\mathbf{m}_i] = p \mathbf{x}_i$$

where p is the probability of retaining a feature.

Also,

$$\text{Cov}(\tilde{\mathbf{x}}_n^T, \tilde{\mathbf{x}}_n) = E[\tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n] - E[\tilde{\mathbf{x}}_n^T] E[\tilde{\mathbf{x}}_n] \quad (7)$$

$$\Rightarrow E[\tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n] = \text{Cov}(\tilde{\mathbf{x}}_n) + \mathbf{p}^2 \mathbf{x}_n^T \mathbf{x}_n \quad (8)$$

Substituting values of $E[\tilde{\mathbf{x}}_n]$, $E[\tilde{\mathbf{x}}_n^T]$ and $E[\tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n]$ in equation (5) gives:

$$E[L(\mathbf{w})] = \sum_{n=1}^N y_n^2 - p y_n \mathbf{w}^T \mathbf{x}_n - p y_n \mathbf{w} \mathbf{x}_n^T + \mathbf{w}^T (\text{Cov}(\tilde{\mathbf{x}}_n) + \mathbf{p}^2 \mathbf{x}_n^T \mathbf{x}_n) \mathbf{w} \quad (9)$$

$$\Rightarrow E[L(\mathbf{w})] = \sum_{n=1}^N (y_n - p \mathbf{w}^T \mathbf{x}_n)^2 + \mathbf{w}^T \text{Cov}(\tilde{\mathbf{x}}_n) \mathbf{w} \quad (10)$$

where $\text{Cov}(\tilde{\mathbf{x}}_n)$ is a diagonal matrix of size $D \times D$ with $p(p-1)\mathbf{x}_{nm}^2$ as diagonal elements, where m ranges from 1 to D .

Let,

$$\sum_{n=1}^N \text{Cov}(\tilde{\mathbf{x}}_n) = \sum$$

Therefore,

$$\mathbb{E}[L(\mathbf{w})] = \sum_{n=1}^N (y_n - p\mathbf{w}^T \mathbf{x}_n)^2 + \mathbf{w}^T \sum \mathbf{w}$$

This equation is similar to the expression of a regularised loss function with $\mathbf{w}^T \sum \mathbf{w}$ as the regularization term. Therefore, minimizing the expected value of $\sum_{n=1}^N (y_n - \mathbf{w}^T \tilde{\mathbf{x}}_n)^2$ is equivalent to minimizing a regularized loss function.

Student Name: Deeksha Arora
 Roll Number: 20111017
 Date: October 30, 2020

For the problem, $\{\hat{\mathbf{B}}, \hat{\mathbf{S}}\} = \arg \min_{\mathbf{B}, \mathbf{S}} \text{TRACE} [(\mathbf{Y} - \mathbf{XBS})^\top (\mathbf{Y} - \mathbf{XBS})]$

The expression for loss function can be written as:

$$L(\mathbf{B}, \mathbf{S}) = \text{TRACE} [(\mathbf{Y} - \mathbf{XBS})^T (\mathbf{Y} - \mathbf{XBS})]$$

$$L(\mathbf{B}, \mathbf{S}) = \text{TRACE} [(\mathbf{Y}^T - \mathbf{S}^T \mathbf{B}^T \mathbf{X}^T) (\mathbf{Y} - \mathbf{XBS})]$$

$$L(\mathbf{B}, \mathbf{S}) = \text{TRACE} [(\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{XBS} - \mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{Y} + \mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{XBS})]$$

$$L(\mathbf{B}, \mathbf{S}) = \text{TRACE} [\mathbf{Y}^T \mathbf{Y}] - \text{TRACE} [\mathbf{Y}^T \mathbf{XBS}] - \text{TRACE} [\mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{Y}] + \text{TRACE} [\mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{XBS}]$$

Keeping \mathbf{S} constant and differentiating $L(\mathbf{B}, \mathbf{S})$ with respect to \mathbf{B} , we get:

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{B}} = 0 - (\mathbf{Y}^T \mathbf{X})^T \mathbf{S}^T - (\mathbf{S} \mathbf{Y}^T \mathbf{X})^T + ((\mathbf{S} \mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{X})^T + (\mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{X})^T \mathbf{S}^T)$$

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{B}} = 0 - \mathbf{X}^T \mathbf{Y} \mathbf{S}^T - \mathbf{X}^T \mathbf{Y} \mathbf{S}^T + (\mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S} \mathbf{S}^T + \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S} \mathbf{S}^T)$$

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{B}} = 2 (\mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S} \mathbf{S}^T) - 2 (\mathbf{X}^T \mathbf{Y} \mathbf{S}^T)$$

Put $\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{B}} = 0$. Therefore,

$$(\mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S} \mathbf{S}^T) = (\mathbf{X}^T \mathbf{Y} \mathbf{S}^T)$$

$$\Rightarrow \mathbf{B} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y} \mathbf{S}^T) (\mathbf{S} \mathbf{S}^T)^{-1}$$

Keeping \mathbf{B} constant and differentiating $L(\mathbf{B}, \mathbf{S})$ with respect to \mathbf{S} , we get:

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{S}} = 0 - \mathbf{B}^T \mathbf{X}^T \mathbf{Y} - (\mathbf{Y}^T \mathbf{X} \mathbf{B})^T + ((\mathbf{S}^T \mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B})^T + (\mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S}))$$

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{S}} = 0 - \mathbf{B}^T \mathbf{X}^T \mathbf{Y} - \mathbf{B}^T \mathbf{X}^T \mathbf{Y} + (\mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S} + \mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S})$$

$$\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{S}} = 2 (\mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S}) - 2 (\mathbf{B}^T \mathbf{X}^T \mathbf{Y})$$

Put $\frac{\partial L(\mathbf{B}, \mathbf{S})}{\partial \mathbf{S}} = 0$. Therefore,

$$(\mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B} \mathbf{S}) = (\mathbf{B}^T \mathbf{X}^T \mathbf{Y})$$

$$\mathbf{S} = (\mathbf{B}^T \mathbf{X}^T \mathbf{X} \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{X}^T \mathbf{Y})$$

$$\Rightarrow \mathbf{S} = ((\mathbf{X} \mathbf{B})^T (\mathbf{X} \mathbf{B}))^{-1} (\mathbf{X} \mathbf{B})^T \mathbf{Y}$$

ALT-OPT algorithm to learn \mathbf{B} and \mathbf{S} for the optimization problem $\{\hat{\mathbf{B}}, \hat{\mathbf{S}}\} = \arg \min_{\mathbf{B}, \mathbf{S}} \text{TRACE} [(\mathbf{Y} - \mathbf{XBS})^\top (\mathbf{Y} - \mathbf{XBS})]$ is as follows:

1. Initialize one the variables. Ex: $\mathbf{S} = \mathbf{S}^{(0)}, t=0$
2. $\mathbf{B}^{(t+1)} = \arg \min_{\mathbf{B}} \mathcal{L}(\mathbf{B}, \mathbf{S}^{(t)})$

$$= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y} \mathbf{S}^{(t)T}) (\mathbf{S}^{(t)} \mathbf{S}^{(t)T})^{-1}$$
3. $\mathbf{S}^{(t+1)} = \arg \min_{\mathbf{S}} \mathcal{L}(\mathbf{B}^{(t+1)}, \mathbf{S})$

$$= \left((\mathbf{X} \mathbf{B}^{(t+1)})^T (\mathbf{X} \mathbf{B}^{(t+1)}) \right)^{-1} (\mathbf{X} \mathbf{B}^{(t+1)})^T \mathbf{Y}$$
4. $t = t + 1$. Go to step 2, if not converged yet.

To compute matrix \mathbf{B} , we need to compute inverse of two matrices, one is $(\mathbf{X}^T \mathbf{X})$ which is a $\mathbf{D} \times \mathbf{D}$ matrix and the other is $(\mathbf{S} \mathbf{S}^T)$ which is $\mathbf{K} \times \mathbf{K}$ matrix. On the other hand, to compute \mathbf{S} , we need to find inverse of $((\mathbf{X} \mathbf{B})^T (\mathbf{X} \mathbf{B}))$ which is a $\mathbf{K} \times \mathbf{K}$ matrix. The computation of inversion of a $\mathbf{K} \times \mathbf{K}$ matrix is common in both the subproblems, but subproblem for \mathbf{B} requires an additional inversion of a $\mathbf{D} \times \mathbf{D}$ matrix which makes computation of matrix \mathbf{B} harder than that of matrix \mathbf{S} . However, the matrix $(\mathbf{X}^T \mathbf{X})$ is same for all iterations of ALT-OPT algorithm. If we assume that inversion of $(\mathbf{X}^T \mathbf{X})$ can be precomputed and cached then we don't need to solve it in every iteration. Hence, if the inversion of $(\mathbf{X}^T \mathbf{X})$ is **cached**, solving the subproblem for \mathbf{B} becomes **equally hard** as solving the subproblem for \mathbf{S} . However, if the inversion of $(\mathbf{X}^T \mathbf{X})$ is **not cached**, then solving the subproblem for \mathbf{B} is **harder** than solving the subproblem for \mathbf{S} .

Student Name: Deeksha Arora

Roll Number: 20111017

Date: October 30, 2020

In Newton's method at each point the quadratic (second-order) approximation of $L(\mathbf{w})$ is minimized. So,

$$\begin{aligned}\mathbf{w}^{(t+1)} &= \arg \min_{\mathbf{w}} \left[L(\mathbf{w}^{(t)}) + \nabla L(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 L(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)}) \right] \\ &= \mathbf{w}^{(t)} - \left(\nabla^2 L(\mathbf{w}^{(t)}) \right)^{-1} \nabla L(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} - \left(\mathbf{H}^{(t)} \right)^{-1} \mathbf{g}^{(t)}\end{aligned}\tag{11}$$

where $\mathbf{g}^{(t)} = \nabla L(\mathbf{w}^{(t)})$ represents the gradient vector for the t^{th} iteration and $\mathbf{H}^{(t)} = \nabla^2 L(\mathbf{w}^{(t)})$ represents the Hessian matrix for the t^{th} iteration.

Ridge regression problem is defined as:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

where $L(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$

Differentiating $L(\mathbf{w})$ with respect to \mathbf{w} we get:

$$\begin{aligned}\nabla L(\mathbf{w}) &= \frac{1}{2} \nabla \left((\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \right) \\ &= \frac{1}{2} (2\mathbf{X}^\top \mathbf{X}\mathbf{w} - 2\mathbf{X}^\top \mathbf{y} + 2\lambda \mathbf{w}) \\ &= \mathbf{X}^\top \mathbf{X}\mathbf{w} - \mathbf{X}^\top \mathbf{y} + \lambda \mathbf{w}\end{aligned}\tag{12}$$

Differentiating $\nabla L(\mathbf{w})$ with respect to \mathbf{w} we get:

$$\begin{aligned}\nabla^2 L(\mathbf{w}) &= \nabla (\mathbf{X}^\top \mathbf{X}\mathbf{w} - \mathbf{X}^\top \mathbf{y} + \lambda \mathbf{w}) \\ &= \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}\end{aligned}\tag{13}$$

Therefore, the gradient vector and the Hessian matrix for the t^{th} iteration are given by:

$$\mathbf{g}^{(t)} = -\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}^{(t)}) + \lambda \mathbf{w}^{(t)}\tag{14}$$

$$\mathbf{H}^{(t)} = \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}\tag{15}$$

Initializing \mathbf{w} as $\mathbf{w}^{(0)}$ and substituting values of $\mathbf{w}^{(t)}$, gradient vector and Hessian matrix in equation (11) gives:

$$\begin{aligned}\Rightarrow \mathbf{w}^{(1)} &= \mathbf{w}^{(0)} - (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \left[-\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X}\mathbf{w}^{(0)} + \lambda \mathbf{w}^{(0)} \right] \\ &= \mathbf{w}^{(0)} - (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \left[-\mathbf{X}^\top \mathbf{y} + (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{w}^{(0)} \right] \\ &= \mathbf{w}^{(0)} + (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} - \mathbf{w}^{(0)} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}\tag{16}$$

The expression obtained for $\mathbf{w}^{(1)}$ is independent of the iteration number t . Therefore, the expression for $\mathbf{w}^{(2)}$ will be same as that of $\mathbf{w}^{(1)}$ i.e.

$$\mathbf{w}^{(2)} = \mathbf{w}^{(1)} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (17)$$

Hence, Newton's method took 2 iterations to converge, where the algorithm is run for the second iteration to identify convergence.

Student Name: Deeksha Arora

Roll Number: 20111017

Date: October 30, 2020

Likelihood:

The given problem is of a categorical distribution. For N dice rolls, $\mathbf{X} = \{x_1, \dots, x_N\}$, where $x_i \in \{1, \dots, 6\}$. So, the Multinoulli distribution is an appropriate choice for the likelihood.

$$p(\mathbf{X} | \boldsymbol{\pi}) = \text{multinoulli}(\mathbf{X} | \boldsymbol{\pi}) \quad (18)$$

If we assume that the data is independent and identically distributed, then likelihood has the form:

$$p(\mathbf{X} | \boldsymbol{\pi}) = \prod_{k=1}^6 \pi_k^{N_k} \quad (19)$$

where $\boldsymbol{\pi}$ is the probability vector such that $\pi_k = p(\mathbf{X} = k)$ and $N_k = \sum_{i=1}^N \mathbb{I}(y_i = k)$ is the number of times event k occurred.

Log-Likelihood:

$$\log(p(\mathbf{X} | \boldsymbol{\pi})) = \sum_{k=1}^6 N_k \log(\pi_k)$$

Prior distribution :

The probability vector $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_6]$ is a K -dimensional vector. So, the Dirichlet distribution is the appropriate prior for this problem, given by:

$$p(\boldsymbol{\pi} | \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^6 \pi_i^{\alpha_i - 1} \quad (20)$$

where the normalizing constant $B(\boldsymbol{\alpha})$ is a multivariate beta function, given by:

$$B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^6 \Gamma(\alpha_i)}{\Gamma\left(\sum_{i=1}^6 \alpha_i\right)}, \quad \boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_6] \quad (21)$$

Log-Prior (without the constant term):

$$\log(p(\boldsymbol{\pi} | \boldsymbol{\alpha})) = \sum_{k=1}^6 (\alpha_k - 1) \log(\pi_k)$$

MAP Estimate:

For derivation of MAP estimate, we need to use a Lagrange multiplier to enforce the equality constraint that $\sum_{k=1}^6 \pi_k = 1$. The constrained objective function or the Lagrangian is given by the log likelihood plus log prior plus the constraint:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\pi}, \lambda) &= \sum_{k=1}^6 N_k \log \pi_k + \sum_{k=1}^6 (\alpha_k - 1) \log \pi_k + \lambda \left(1 - \sum_{k=1}^6 \pi_k\right) \\ &= \sum_{k=1}^6 (N_k + \alpha_k - 1) \log \pi_k + \lambda \left(1 - \sum_{k=1}^6 \pi_k\right) \end{aligned} \quad (22)$$

Taking derivative of $\mathcal{L}(\boldsymbol{\pi}, \lambda)$ with respect to λ gives:

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \left(1 - \sum_{k=1}^6 \pi_k\right) = 0 \quad (23)$$

Taking derivative of $\mathcal{L}(\boldsymbol{\pi}, \lambda)$ with respect to π_k gives:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \pi_k} &= \frac{N_k + \alpha_k - 1}{\pi_k} - \lambda = 0 \\ \Rightarrow N_k + \alpha_k - 1 &= \lambda \pi_k \end{aligned} \quad (24)$$

Using the sum-to-one constraint to solve for λ :

$$\begin{aligned} \sum_k N_k + \alpha_k - 1 &= \lambda \sum_k \pi_k \\ \Rightarrow \lambda &= N + \alpha_0 - 6 \end{aligned} \quad (25)$$

where $\alpha_0 \triangleq \sum_{k=1}^6 \alpha_k$ is the equivalent sample size of the prior. Substituting value of λ in equation (24) gives:

$$\hat{\pi}_k = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - 6} \quad (26)$$

which is the required MAP estimate.

MAP solution is expected to be better than MLE solution in the following cases:

1. When number of training examples is less. In this case, the MLE solution tends to overfit and is often unreliable. On the other hand, MAP solution will give more accurate results because it uses the regularizer term which prevents overfitting.
2. When the training data contains examples from a specific class. For example, considering the dice roll problem where in every roll, the dice shows face 6. In this case, the MLE estimate will give $\pi_6 = 1$, and the probability for all other faces is 0. MAP solution would give more reliable estimation in this case.

Full posterior distribution:

The posterior distribution is given by:

$$\begin{aligned} p(\boldsymbol{\pi} \mid \mathbf{X}, \boldsymbol{\alpha}) &= \frac{p(\mathbf{X} \mid \boldsymbol{\pi}) * p(\boldsymbol{\pi} \mid \boldsymbol{\alpha})}{p(\mathbf{X})} \\ &= \frac{\prod_{k=1}^6 \pi_k^{N_k} * \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^6 \pi_i^{\alpha_i - 1}}{\int (\prod_{k=1}^6 \pi_k^{N_k} * \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^6 \pi_i^{\alpha_i - 1}) d\boldsymbol{\pi}} \end{aligned} \quad (27)$$

Since the denominator is constant, the posterior can be represented as:

$$\begin{aligned} p(\boldsymbol{\pi} \mid \mathbf{X}, \boldsymbol{\alpha}) &\propto p(\mathbf{X} \mid \boldsymbol{\pi}) p(\boldsymbol{\pi} \mid \boldsymbol{\alpha}) \\ &\propto \prod_{k=1}^6 \pi_k^{N_k} \pi_k^{\alpha_k - 1} \\ &\propto \prod_{k=1}^6 \pi_k^{\alpha_k + N_k - 1} \\ &= \text{Dirichlet}(\boldsymbol{\pi} \mid \alpha_1 + N_1, \dots, \alpha_6 + N_6) \end{aligned} \quad (28)$$

Yes, given the above posterior we can get the MLE and MAP estimate without solving the MLE and MAP optimization problems.

MAP estimate from the posterior:

The MAP estimate will be the mode of the above Dirichlet distribution and is given by:

$$\hat{\pi}_k(MAP) = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - 6} \quad (29)$$

MLE estimate from the posterior:

When the prior is uniform, MAP estimate is equivalent to MLE estimate. Therefore, to get the MLE solution, we need to make the prior uniform. This can be done by equating each component of α to 1 and equation (29) will give the MLE solution as:

$$\hat{\pi}_k(MLE) = \frac{N_k}{N} \quad (30)$$