

EECE7397 – Homework 1

1.

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_j(x_n)^j$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left[y(x_n, \mathbf{w}) - t_n \right]^2$$

$$\Rightarrow E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right]^2$$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_{n=1}^N 2 \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] \frac{\partial}{\partial w_i} \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right]$$

$$= \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] (x_n)^i \left(\text{As } \frac{\partial}{\partial w_i} w_j(x_n)^j = \delta_{ij} (x_n)^i \text{ if } j = i \right)$$

$$\frac{\partial E}{\partial w_i} = 0$$

$$\Rightarrow \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] (x_n)^i = 0$$

$$\Rightarrow \sum_{n=1}^N \sum_{j=0}^M w_j(x_n)^j (x_n)^i - \sum_{n=1}^N t_n (x_n)^i = 0$$

$$\Rightarrow \sum_{n=1}^N \sum_{j=0}^M w_j(x_n)^{j+i} - \sum_{n=1}^N t_n (x_n)^i = 0$$

$$\Rightarrow \sum_{j=0}^M w_j \sum_{n=1}^N (x_n)^{j+i} - \sum_{n=1}^N t_n (x_n)^i = 0$$

$$\Rightarrow \sum_{j=0}^M w_j \sum_{n=1}^N (x_n)^{j+i} = \sum_{n=1}^N t_n (x_n)^i$$

$$\Rightarrow \sum_{j=0}^M A_{ij} w_j = T_i \left(A_{ij} = \sum_{n=1}^N (x_n)^{j+i}, \quad T_i = \sum_{n=1}^N t_n (x_n)^i \right)$$

In matrix notation:

$$\mathbf{A} \mathbf{w} = \mathbf{T}$$

$$\Rightarrow \mathbf{w} = \mathbf{A}^{-1} \mathbf{T}$$

2.

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_j(x_n)^j$$

$$\begin{aligned}\tilde{E}(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \left[y(x_n, \mathbf{w}) - t_n \right]^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \\ \Rightarrow \tilde{E}(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right]^2 + \frac{\lambda}{2} \sum_{j=0}^M w_j^2\end{aligned}$$

$$\begin{aligned}\frac{\partial \tilde{E}}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right]^2 + \frac{\partial}{\partial w_i} \frac{\lambda}{2} \sum_{j=0}^M w_j^2 \\ &= \frac{1}{2} \sum_{n=1}^N 2 \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] \frac{\partial}{\partial w_i} \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] + \frac{\lambda}{2} \frac{\partial}{\partial w_i} \sum_{j=0}^M w_j^2 \\ &= \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] (x_n)^i + \lambda w_i \left(\text{As } \frac{\partial}{\partial w_i} w_j(x_n)^j = \delta_{ij}(x_n)^i \text{ if } j = i \right)\end{aligned}$$

$$\begin{aligned}\frac{\partial \tilde{E}}{\partial w_i} &= 0 \\ \Rightarrow \sum_{n=1}^N \left[\sum_{j=0}^M w_j(x_n)^j - t_n \right] (x_n)^i + \lambda w_i &= 0 \\ \Rightarrow \sum_{n=1}^N \sum_{j=0}^M w_j(x_n)^j (x_n)^i - \sum_{n=1}^N (x_n)^i t_n + \lambda w_i &= 0 \\ \Rightarrow \sum_{n=1}^N \sum_{j=0}^M w_j(x_n)^{j+i} - \sum_{n=1}^N (x_n)^i t_n + \lambda w_i &= 0 \\ \Rightarrow \sum_{j=0}^M w_j \sum_{n=1}^N (x_n)^{j+i} - \sum_{n=1}^N (x_n)^i t_n + \lambda w_i &= 0 \\ \Rightarrow \sum_{j=0}^M w_j \sum_{n=1}^N (x_n)^{j+i} + \lambda w_i &= \sum_{n=1}^N (x_n)^i t_n \\ \Rightarrow \sum_{j=0}^M A_{ij} w_j + \lambda w_i &= T_i \left(A_{ij} = \sum_{n=1}^N (x_n)^{j+i}, \quad T_i = \sum_{n=1}^N t_n (x_n)^i \right)\end{aligned}$$

In matrix notation:

$$\begin{aligned}\mathbf{A} \mathbf{w} + \lambda \mathbf{w} &= \mathbf{T} \\ \Rightarrow (\mathbf{A} + \lambda \mathbf{I}) \mathbf{w} &= \mathbf{T} \\ \Rightarrow \mathbf{w} &= (\mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{T}\end{aligned}$$

The regularization penalty term shrinks the parameter vector \mathbf{w} toward 0 or smaller magnitudes, preventing large parameter values that can lead to overfitting. Another example of regularization is lasso, which adds a $\lambda \|\mathbf{w}\|$ penalty term. This is similar to the ridge regularization term $\frac{\lambda}{2} \|\mathbf{w}\|^2$ in that both push parameters toward zero; but, the lasso regularization penalty promotes sparsity, resulting in some coefficients resulting in 0.

3.

Assuming for a simple function $f(x)$, change of variable using non-linear function $g(y)$ results in $h(y) = f(g(y))$.

$$\frac{d}{dy} [f(g(y))] = f'(g(y)) g'(y)$$

$h(y)$ has maximum at \hat{y} (mode \hat{y}):

$$h'(\hat{y}) = f'(g(\hat{y})) g'(\hat{y}) = 0$$

Assuming $g'(\hat{y}) \neq 0$:

$$f'(g(\hat{y})) = 0$$

$f(x)$ has maximum at \hat{x} (mode \hat{x}):

$$f'(\hat{x}) = 0$$

As $f'(\hat{x}) = 0$ and $f'(g(\hat{y})) = 0$:

$$\hat{x} = g(\hat{y})$$

The \hat{y} that maximizes h in y maps directly to the \hat{x} that maximizes f in x via $\hat{x} = g(\hat{y})$, which is the maximum of f in x .

For a probability density function $p_x(x)$ change of variable using non-linear function $g(y)$ results in $p_y(y) = p_x(g(y)) |g'(y)|$

$$\begin{aligned} \frac{d}{dy} p_y(y) &= \frac{d}{dy} [p_x(g(y)) |g'(y)|] \\ &= \frac{d}{dy} [p_x(g(y))] \cdot |g'(y)| + p_x(g(y)) \frac{d}{dy} [|g'(y)|] \\ &= p'_x(g(y)) g'(y) \cdot |g'(y)| + p_x(g(y)) \cdot \text{sgn}[g'(y)] g''(y) \end{aligned}$$

$p_x(x)$ has a maximum at \hat{x} (mode \hat{x}):

$$p'_x(\hat{x}) = 0$$

$p_y(y)$ has a maximum at \hat{y} (mode \hat{y}):

$$\begin{aligned} p'_y(\hat{y}) &= 0 \\ \Rightarrow p'_x(g(\hat{y})) g'(\hat{y}) \cdot |g'(\hat{y})| + p_x(g(\hat{y})) \cdot \text{sgn}[g'(\hat{y})] g''(\hat{y}) &= 0 \\ \Rightarrow p'_x(g(\hat{y})) g'(\hat{y}) \cdot g'(\hat{y}) \cdot \text{sgn}[g'(\hat{y})] + p_x(g(\hat{y})) \text{sgn}[g'(\hat{y})] g''(\hat{y}) &= 0 \\ \Rightarrow \text{sgn}[g'(\hat{y})] \left(p'_x(g(\hat{y})) g'(\hat{y})^2 + p_x(g(\hat{y})) g''(\hat{y}) \right) &= 0 \\ \Rightarrow p'_x(g(\hat{y})) g'(\hat{y})^2 + p_x(g(\hat{y})) g''(\hat{y}) &= 0 \\ \Rightarrow p'_x(g(\hat{y})) g'(\hat{y})^2 = -p_x(g(\hat{y})) g''(\hat{y}) \end{aligned}$$

Because $g''(\hat{y}) \neq 0$, the expression does not reduce to $p'_x(g(\hat{y})) = 0$. The extra term $p_x(g(\hat{y})) g''(\hat{y})$, arising from the derivative of the Jacobian factor $|g'(y)|$, shifts the maximum away from the point where $p'_x(x) = 0$. The mode of $p_y(y)$ does not necessarily map to the mode of $p_x(x)$ via $x = g(y)$.

For a linear transformation, assume function $g(y) = a y + b$,

$$g'(y) = a, \quad g''(y) = 0.$$

For a probability density function $p_x(x)$ change of variable using linear function $g(y)$ results in $p_y(y) = p_x(a y + b) |a|$.

$$\begin{aligned}\frac{d}{dy} p_y(y) &= \frac{d}{dy} [p_x(a y + b) |a|] \\ &= |a| \frac{d}{dy} [p_x(a y + b)] \\ &= |a| [p'_x(a y + b) a] \\ &= a |a| p'_x(a y + b).\end{aligned}$$

$p_x(x)$ has a maximum at \hat{x} (mode \hat{x}):

$$p'_x(\hat{x}) = 0$$

$p_y(y)$ has a maximum at \hat{y} (mode \hat{y}):

$$\begin{aligned}p'_y(\hat{y}) &= 0 \\ \Rightarrow a |a| p'_x(a \hat{y} + b) &= 0 \\ \Rightarrow p'_x(a \hat{y} + b) &= 0\end{aligned}$$

As $p'_x(\hat{x}) = 0$ and $p'_x(a \hat{y} + b) = 0$:

$$\hat{x} = a \hat{y} + b$$

For a general nonlinear $g(y)$, the transformed density includes the Jacobian factor $|g'(y)|$ which results in When you an extra term involving $g''(y)$ shifting the mode when finding the location of the maximum. Consequently, the value \hat{y} that maximizes $p_y(y)$ does not, in general, map via $\hat{x} = g(\hat{y})$ to the value \hat{x} that maximizes $p_x(x)$. This shift in the mode occurs because the Jacobian term $|g'(y)|$ affects where the maximum occurs.

If $g(y) = a y + b$ with $a \neq 0$, then $g'(y) = a$ and $g''(y) = 0$, so the Jacobian factor $|g'(y)| = |a|$ is a constant. Thus, maximizing $p_x(a y + b) |a|$ is equivalent to maximizing $p_x(a y + b)$. The mode \hat{y} of $p_y(y)$ and the mode \hat{x} of $p_x(x)$ map via $\hat{x} = a \hat{y} + b$.

4.

$$E[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

$$\|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 = (\mathbf{y}(\mathbf{x}) - \mathbf{t})^\top (\mathbf{y}(\mathbf{x}) - \mathbf{t}) = \sum_{i=1}^M [y_i(\mathbf{x}) - t_i]^2$$

$$\begin{aligned}E[L] &= \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \\ &= \iint \sum_{i=1}^M [y_i(\mathbf{x}) - t_i]^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}\end{aligned}$$

$$\begin{aligned}\delta E[L] &= \frac{\delta}{\delta y_j(\mathbf{x})} \iint \sum_{i=1}^M [y_i(\mathbf{x}) - t_i]^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \\ &= \iint \frac{\delta}{\delta y_j(\mathbf{x})} \left[\sum_{i=1}^M (y_i(\mathbf{x}) - t_i)^2 \right] p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \\ &= \iint 2 [y_j(\mathbf{x}) - t_j] p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \\ &= \int 2 [y_j(\mathbf{x}) - t_j] p(\mathbf{x}, \mathbf{t}) d\mathbf{t}\end{aligned}$$

$$\begin{aligned}
& \int 2 [y_j(\mathbf{x}) - t_j] p(\mathbf{x}, \mathbf{t}) d\mathbf{t} = 0 \\
\Rightarrow y_j(\mathbf{x}) \int p(\mathbf{x}, \mathbf{t}) d\mathbf{t} - \int t_j p(\mathbf{x}, \mathbf{t}) d\mathbf{t} &= 0 \\
\Rightarrow y_j(\mathbf{x}) \int p(\mathbf{x}, \mathbf{t}) d\mathbf{t} &= \int t_j p(\mathbf{x}, \mathbf{t}) d\mathbf{t} \\
\Rightarrow y_j(\mathbf{x}) &= \frac{\int t_j p(\mathbf{x}, \mathbf{t}) d\mathbf{t}}{\int p(\mathbf{x}, \mathbf{t}) d\mathbf{t}} \\
\Rightarrow y_j(\mathbf{x}) &= \frac{\int t_j p(\mathbf{x}, \mathbf{t}) d\mathbf{t}}{p(\mathbf{x})} \quad (\text{As } p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{t}) d\mathbf{t}) \\
\Rightarrow y_j(\mathbf{x}) &= \frac{\int t_j p(\mathbf{t} | \mathbf{x}) p(\mathbf{x}) d\mathbf{t}}{p(\mathbf{x})} \quad (\text{As } p(\mathbf{x}, \mathbf{t}) = p(\mathbf{t} | \mathbf{x}) p(\mathbf{x})) \\
\Rightarrow y_j(\mathbf{x}) &= \int t_j p(\mathbf{t} | \mathbf{x}) d\mathbf{t} \\
\Rightarrow y_j(\mathbf{x}) &= \mathbf{E}[t_j | \mathbf{x}] \quad (\text{As } \mathbf{E}[x | y] = \int x p(x | y) dx)
\end{aligned}$$

For multiple target variables $\mathbf{t} = (t_1, \dots, t_M)$:

$$\mathbf{y}(\mathbf{x}) = (y_1(\mathbf{x}), \dots, y_M(\mathbf{x})) = \mathbf{E}[\mathbf{t} | \mathbf{x}],$$

$$y_j(\mathbf{x}) = \int t_j p(\mathbf{t} | \mathbf{x}) d\mathbf{t} = \mathbf{E}[t_j | \mathbf{x}]$$

For a single target variable t ($j = 1$ as $M = 1$) :

$$y(\mathbf{x}) = \int t p(t | \mathbf{x}) dt = \mathbf{E}[t | \mathbf{x}]$$

5.

$$\binom{N}{m} = \frac{N!}{m!(N-m)!}$$

$$\binom{N}{m-1} = \frac{N!}{(m-1)![N-(m-1)]!} = \frac{N!}{(m-1)!(N-m+1)!}$$

$$\begin{aligned}
\binom{N}{m} + \binom{N}{m-1} &= \frac{N!}{m!(N-m)!} + \frac{N!}{(m-1)!(N-m+1)!} \\
&= \frac{N!}{m(m-1)!(N-m)!} + \frac{N!}{(m-1)!(N-m+1)(N-m)!} \\
&\quad (\text{As } m! = m(m-1)! \text{ and } (N-m+1)! = (N-m+1)(N-m)!) \\
&= \frac{N!}{(N-m)!(m-1)!} \frac{1}{m} + \frac{N!}{(N-m)!(m-1)!} \frac{1}{N-m+1} \\
&= \frac{N!}{(N-m)!(m-1)!} \left(\frac{1}{m} + \frac{1}{N-m+1} \right) \\
&= \frac{N!}{(N-m)!(m-1)!} \left(\frac{N-m+1+m}{m(N-m+1)} \right) \\
&= \frac{N!}{(N-m)!(m-1)!} \left(\frac{N+1}{m(N-m+1)} \right).
\end{aligned}$$

$$\begin{aligned}
\binom{N+1}{m} &= \frac{(N+1)!}{m! [(N+1)-m]!} = \frac{(N+1)!}{m! (N+1-m)!} \\
&= \frac{(N+1) N!}{m (m-1)! (N+1-m) (N-m)!} \\
&= \frac{N!}{(N-m)! (m-1)!} \frac{N+1}{m (N+1-m)}.
\end{aligned}$$

As the expressions for $\binom{N}{m} + \binom{N}{m-1}$ and $\binom{N+1}{m}$ match, identity is proven:

$$\binom{N}{m} + \binom{N}{m-1} = \binom{N+1}{m}$$

Given the identity:

$$\binom{N}{m} + \binom{N}{m-1} = \binom{N+1}{m}$$

Proof of the Binomial Theorem by mathematical induction:

$$(1+x)^N = \sum_{m=0}^N \binom{N}{m} x^m$$

Base Case N=0:

$$\begin{aligned}
(1+x)^N &= \sum_{m=0}^N \binom{N}{m} x^m \\
\Rightarrow (1+x)^0 &= \sum_{m=0}^0 \binom{0}{m} x^m \\
\Rightarrow (1+x)^0 &= \sum_{m=0}^0 \binom{0}{0} x^0 \\
\Rightarrow 1 &= 1
\end{aligned}$$

The base case holds for N=0.

Inductive Hypothesis:

Assume the theorem holds for some $N \geq 0$, then it has to proven in the inductive step that it holds for $N+1$

$$(1+x)^N = \sum_{m=0}^N \binom{N}{m} x^m$$

Inductive Step:

To prove it holds for $N+1$:

$$(1+x)^{N+1} = \sum_{m=0}^{N+1} \binom{N+1}{m} x^m$$

$$\begin{aligned}
& (1+x)^{N+1} \\
&= (1+x)^N \cdot (1+x) \\
&= \sum_{m=0}^N \binom{N}{m} x^m \cdot (1+x) \\
&= \sum_{m=0}^N \binom{N}{m} x^m + \sum_{m=0}^N \binom{N}{m} x^m \cdot x \\
&= \sum_{m=0}^N \binom{N}{m} x^m + \sum_{m=0}^N \binom{N}{m} x^{m+1}
\end{aligned}$$

Term x^k for $1 \leq k \leq N$:

In the first sum, the term $\binom{N}{m} x^m$ matches x^k when $m = k$, giving the coefficient $\binom{N}{k}$. In the second sum, $\binom{N}{m} x^{m+1}$ matches x^k when $m + 1 = k$, giving the coefficient $\binom{N}{k-1}$. By the identity $\binom{N}{k} + \binom{N}{k-1} = \binom{N+1}{k}$ the total coefficient of x^k is $\binom{N+1}{k}$. The constant term x^0 (when $m = 0$) comes from the first sum $\binom{N}{0} x^0 = 1$, and the x^{N+1} term comes from the second sum $\binom{N}{N} x^{N+1} = x^{N+1}$ when $m = N$. The total coefficient of the powers x^1, x^2, \dots, x^N . Thus, the expression becomes:

$$(1+x)^{N+1} = \sum_{m=0}^{N+1} \binom{N+1}{m} x^m$$

$$\begin{aligned}
& \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= (1-\mu)^N \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{-m} \\
&= (1-\mu)^N \sum_{m=0}^N \binom{N}{m} \left(\frac{\mu}{1-\mu}\right)^m
\end{aligned}$$

Let x be $\frac{\mu}{1-\mu}$.

$$(1+x)^N = \left(1 + \frac{\mu}{1-\mu}\right)^N = \left(\frac{1-\mu+\mu}{1-\mu}\right)^N = \left(\frac{1}{1-\mu}\right)^N = (1-\mu)^{-N}$$

From the binomial theorem:

$$(1+x)^N = \sum_{m=0}^N \binom{N}{m} x^m$$

Thus, $\sum_{m=0}^N \binom{N}{m} x^m = (1-\mu)^{-N}$

$$\begin{aligned}
& (1-\mu)^N \sum_{m=0}^N \binom{N}{m} \left(\frac{\mu}{1-\mu}\right)^m \\
&= (1-\mu)^N \sum_{m=0}^N \binom{N}{m} x^m \\
&= (1-\mu)^N (1-\mu)^{-N} \\
&= 1
\end{aligned}$$

Thus, $\sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} = 1$ is proven.

$$\begin{aligned}
\mathbf{E}[\mathbf{X}] &= \sum_{m=0}^N m \cdot \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N m \cdot \frac{N!}{(N-m)!m!} \cdot \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N m \cdot \frac{N(N-1)!}{(N-m)! \cdot m(m-1)!} \cdot \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N \frac{N(N-1)!}{(N-m)!(m-1)!} \cdot \mu^m (1-\mu)^{N-m} \\
&= N \sum_{m=1}^N \frac{(N-1)!}{(N-m)!(m-1)!} \cdot \mu^m (1-\mu)^{N-m}
\end{aligned}$$

Let $k = m - 1$, then $m = k + 1$. As m goes from 1 to N , k goes from 0 to $N - 1$.

$$\begin{aligned}
&N \sum_{m=1}^N \frac{(N-1)!}{(N-m)!(m-1)!} \cdot \mu^m (1-\mu)^{N-m} \\
&= N \sum_{k=0}^{N-1} \frac{(N-1)!}{(N-(k+1))!((k+1)-1)!} \cdot \mu^{k+1} (1-\mu)^{N-(k+1)} \\
&= N \sum_{k=0}^{N-1} \frac{(N-1)!}{(N-1-k)!k!} \cdot \mu^{k+1} (1-\mu)^{(N-1)-k} \\
&= N \sum_{k=0}^{N-1} \frac{(N-1)!}{(N-1-k)!k!} \cdot \mu^k \mu (1-\mu)^{(N-1)-k} \\
&= N \mu \sum_{k=0}^{N-1} \frac{(N-1)!}{(N-1-k)!k!} \cdot \mu^k (1-\mu)^{(N-1)-k} \\
&= N \mu \sum_{k=0}^{N-1} \binom{N-1}{k} \mu^k (1-\mu)^{(N-1)-k} \\
&= N \mu
\end{aligned}$$

$$\left(As(1+x)^N = \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} = 1 \right)$$

$$\left(As(\mu + (1-\mu))^{N-1} = \sum_{k=0}^{N-1} \binom{N-1}{k} \mu^k (1-\mu)^{(N-1)-k} = 1 \right)$$

$$\begin{aligned}
\mathbf{E}[\mathbf{X}^2] &= \sum_{m=0}^N m^2 \cdot \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N (m(m-1) + m) \cdot \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \left(\sum_{m=0}^N m(m-1) + \sum_{m=0}^N m \right) \cdot \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N m(m-1) \binom{N}{m} \mu^m (1-\mu)^{N-m} + \sum_{m=0}^N m \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N m(m-1) \cdot \frac{N!}{(N-m)!m!} \mu^m (1-\mu)^{N-m} + \sum_{m=0}^N m \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N m(m-1) \cdot \frac{N(N-1)(N-2)!}{(N-m)!m(m-1)(m-2)!} \mu^m (1-\mu)^{N-m} + \sum_{m=0}^N m \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=0}^N \frac{N(N-1)(N-2)!}{(N-m)!(m-2)!} \mu^m (1-\mu)^{N-m} + \sum_{m=0}^N m \binom{N}{m} \mu^m (1-\mu)^{N-m} \\
&= N(N-1) \sum_{m=0}^N \frac{(N-2)!}{(N-m)!(m-2)!} \mu^m (1-\mu)^{N-m} + \sum_{m=0}^N m \binom{N}{m} \mu^m (1-\mu)^{N-m}
\end{aligned}$$

Let $k = m - 2$, then $m = k + 2$. As m goes from 2 to N , k goes from 0 to $N - 2$.

$$\begin{aligned}
&N(N-1) \sum_{m=0}^N \frac{(N-2)!}{(N-m)!(m-2)!} \mu^m (1-\mu)^{N-m} \\
&= N(N-1) \sum_{k=0}^{N-2} \frac{(N-2)!}{(N-(k+2))!(k+2-2)!} \mu^{k+2} (1-\mu)^{N-(k+2)} \\
&= N(N-1) \sum_{k=0}^{N-2} \frac{(N-2)!}{(N-2-k)!k!} \mu^{k+2} (1-\mu)^{N-(k+2)} \\
&= N(N-1) \sum_{k=0}^{N-2} \binom{N-2}{k} \mu^{k+2} (1-\mu)^{N-(k+2)} \\
&= N(N-1) \sum_{k=0}^{N-2} \binom{N-2}{k} \mu^k \mu^2 (1-\mu)^{N-(k+2)} \\
&= N(N-1) \mu^2 \sum_{k=0}^{N-2} \binom{N-2}{k} \mu^k (1-\mu)^{(N-2)-k} \\
&= N(N-1) \mu^2 \left(As(1+x)^N = \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} = 1 \right) \\
&\left(As(\mu + (1-\mu))^{N-2} = \sum_{k=0}^{N-2} \binom{N-2}{k} \mu^k (1-\mu)^{(N-2)-k} = 1 \right)
\end{aligned}$$

Thus,

$$\mathbf{E}[\mathbf{X}^2] = N(N-1)\mu^2 + N\mu$$

$$\begin{aligned}
\text{var}(\mathbf{X}) &= \mathbf{E}[\mathbf{X}^2] - \mathbf{E}[\mathbf{X}]^2 \\
&= N(N-1)\mu^2 + N\mu - (N\mu)^2 \\
&= N^2\mu^2 - N\mu^2 + N\mu - N^2\mu^2 \\
&= N\mu - N\mu^2 \\
&= N\mu(1-\mu)
\end{aligned}$$

6.

$$\text{Beta}(\mu \mid a, b) = p(\mu) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

$$\begin{aligned}
E[\mu] &= \int_0^1 \mu p(\mu) d\mu \\
&= \int_0^1 \mu \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \mu \mu^{a-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \mu^{(a+1)-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+1)\Gamma(b)}{\Gamma((a+1)+b)} \quad (\text{As } \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}) \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{a\Gamma(a)\Gamma(b)}{(a+b)\Gamma(a+b)} \quad (\text{As } \Gamma(a+1) = a\Gamma(a), \Gamma(a+b+1) = (a+b)\Gamma(a+b)) \\
&= \frac{a}{a+b}
\end{aligned}$$

$$\begin{aligned}
E[\mu^2] &= \int_0^1 \mu^2 p(\mu) d\mu \\
&= \int_0^1 \mu^2 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \mu^2 \mu^{a-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \mu^{(a+2)-1} (1-\mu)^{b-1} d\mu \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+2)\Gamma(b)}{\Gamma((a+2)+b)} \quad (\text{As } \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}) \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} \\
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{(a+1)a\Gamma(a)\Gamma(b)}{(a+b+1)(a+b)\Gamma(a+b)} \\
&(\text{As } \Gamma((a+1)+1) = (a+1)\Gamma(a+1) = (a+1)a\Gamma(a)) \\
&(\text{As } \Gamma((a+b+1)+1) = (a+b+1)(a+b)\Gamma(a+b)) \\
&= \frac{a(a+1)}{(a+b)(a+b+1)}
\end{aligned}$$

$$\begin{aligned}
\text{Var}[\mu] &= E[\mu^2] - (E[\mu])^2 \\
&= \frac{a(a+1)}{(a+b)(a+b+1)} - \left(\frac{a}{a+b}\right)^2 \\
&= \frac{a(a+1)}{(a+b)(a+b+1)} - \frac{a^2}{(a+b)^2} \\
&= \frac{a(a+1)}{(a+b)(a+b+1)} \cdot \frac{(a+b)}{(a+b)} - \frac{a^2}{(a+b)^2} \cdot \frac{(a+b+1)}{(a+b+1)} \\
&= \frac{a(a+1)(a+b) - a^2(a+b+1)}{(a+b)^2(a+b+1)} \\
&= \frac{(a^3 + a^2b + a^2 + ab) - (a^3 + a^2b + a^2)}{(a+b)^2(a+b+1)} \\
&= \frac{ab}{(a+b)^2(a+b+1)}
\end{aligned}$$

$$p(\mu) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$f(\mu) = \mu^{a-1} (1-\mu)^{b-1}$$

As $\frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ is a constant (independent of μ), the value of μ that maximizes $p(\mu)$ also maximizes $f(\mu)$. To find the mode, the constant is ignored and $f(\mu)$ is maximized.

$$\begin{aligned}
& \frac{d}{d\mu} [\mu^{a-1}(1-\mu)^{b-1}] \\
&= \mu^{a-1} \frac{d}{d\mu} [(1-\mu)^{b-1}] + (1-\mu)^{b-1} \frac{d}{d\mu} [\mu^{a-1}] \\
&= \mu^{a-1} [(b-1)(1-\mu)^{b-2}(-1)] + (1-\mu)^{b-1} [(a-1)\mu^{a-2}] \\
&= -\mu^{a-1}(b-1)(1-\mu)^{b-2} + \mu^{a-2}(1-\mu)^{b-1}(a-1) \\
&= \mu^{a-2}(1-\mu)^{b-2} [-\mu(b-1) + (1-\mu)(a-1)] \\
&= \mu^{a-2}(1-\mu)^{b-2} [(a-1) - (a+b-2)\mu]
\end{aligned}$$

$$\begin{aligned}
& \frac{df}{d\mu} = 0 \\
& \Rightarrow \mu^{a-2}(1-\mu)^{b-2} [(a-1) - (a+b-2)\mu] = 0 \\
& \Rightarrow [(a-1) - (a+b-2)\mu] = 0 \\
& \Rightarrow \mu = \frac{a-1}{a+b-2}
\end{aligned}$$

$$mode[\mu] = \frac{a-1}{a+b-2}$$

7.

$$\begin{aligned}
\mathbf{E}_{x,y}[x+ay] &= \iint (x+ay) p(x,y) dx dy \\
&= \iint (x+ay) p(x) p(y) dx dy \\
&= \iint x p(x) p(y) dx dy + \iint a y p(x) p(y) dx dy \\
&= \left(\int x p(x) dx \right) \left(\int p(y) dy \right) + a \left(\int p(x) dx \right) \left(\int y p(y) dy \right) \\
&= \int x p(x) dx + a \int y p(y) dy \quad (\text{As } \int p(x) dx = 1 \text{ and } \int p(y) dy = 1) \\
&= \mathbf{E}_x[x] + a \mathbf{E}_y[y] \quad (\text{As } \mathbf{E}_x[x] = \int x p(x) dx, \mathbf{E}_y[y] = \int y p(y) dy)
\end{aligned}$$

$$\begin{aligned}
\text{var}_{x,y}[x + ay] &= \mathbf{E}_{x,y}[(x + ay)^2] - \mathbf{E}_{x,y}[x + ay]^2 \\
&= \iint (x + ay)^2 p(x, y) \, dx \, dy - \left(\mathbf{E}_x[x] + a \mathbf{E}_y[y] \right)^2 \\
&= \iint (x^2 + 2axy + a^2 y^2) p(x)p(y) \, dx \, dy - \left(\mathbf{E}_x[x] + a \mathbf{E}_y[y] \right)^2 \\
&= \left(\int x^2 p(x) \, dx \right) \left(\int p(y) \, dy \right) + 2a \left(\int x p(x) \, dx \right) \left(\int y p(y) \, dy \right) \\
&\quad + a^2 \left(\int p(x) \, dx \right) \left(\int y^2 p(y) \, dy \right) - \left(\mathbf{E}_x[x] + a \mathbf{E}_y[y] \right)^2 \\
&= \mathbf{E}_x[x^2] + 2a \mathbf{E}_x[x] \mathbf{E}_y[y] + a^2 \mathbf{E}_y[y^2] - \left(\mathbf{E}_x[x] + a \mathbf{E}_y[y] \right)^2 \\
&\quad \left(\text{As } \int p(x) \, dx = 1, \int p(y) \, dy = 1, \int x p(x) \, dx = \mathbf{E}_x[x], \int y p(y) \, dy = \mathbf{E}_y[y] \right) \\
&= \mathbf{E}_x[x^2] + 2a \mathbf{E}_x[x] \mathbf{E}_y[y] + a^2 \mathbf{E}_y[y^2] - \left(\mathbf{E}_x[x]^2 + 2a \mathbf{E}_x[x] \mathbf{E}_y[y] + a^2 \mathbf{E}_y[y]^2 \right) \\
&= (\mathbf{E}_x[x^2] - \mathbf{E}_x[x]^2) + a^2 (\mathbf{E}_y[y^2] - \mathbf{E}_y[y]^2) \\
&= \text{var}_x[x] + a^2 \text{var}_y[y] \\
&\quad \left(\text{As } \text{var}_x[x] = \mathbf{E}_x[x^2] - (\mathbf{E}_x[x])^2 \right)
\end{aligned}$$

$$\begin{aligned}
\mathbf{E}[x] &= \iint x p(x, y) \, dx \, dy \\
&= \iint x p(x | y) p(y) \, dx \, dy \quad (\text{As } p(x, y) = p(x | y) p(y)) \\
&= \int \left(\int x p(x | y) \, dx \right) p(y) \, dy \\
&= \int \mathbf{E}_x[x | y] p(y) \, dy \quad (\text{As } \mathbf{E}_x[x] = \int x p(x) \, dx) \\
&= \mathbf{E}_y[\mathbf{E}_x[x | y]] \quad (\text{As } \mathbf{E}_y[y] = \int y p(y) \, dx)
\end{aligned}$$

$$\begin{aligned}
\text{var}_x[x] &= \mathbf{E}_x[x^2] - (\mathbf{E}_x[x])^2 \\
&= \mathbf{E}_y[\mathbf{E}_x[x^2 | y]] - (\mathbf{E}_y[\mathbf{E}_x[x | y]])^2 \quad (\text{As } \mathbf{E}[x] = \mathbf{E}_y[\mathbf{E}_x[x | y]]) \\
&= \mathbf{E}_y[\text{var}_x[x | y] + \mathbf{E}_x[x | y]^2] - (\mathbf{E}_y[\mathbf{E}_x[x | y]])^2 \quad (\text{As } \mathbf{E}_x[x^2] = \text{var}_x[x] + (\mathbf{E}_x[x])^2) \\
&= \mathbf{E}_y[\text{var}_x[x | y]] + \mathbf{E}_y[\mathbf{E}_x[x | y]^2] - (\mathbf{E}_y[\mathbf{E}_x[x | y]])^2 \\
&= \mathbf{E}_y[\text{var}_x[x | y]] + \text{var}_y[\mathbf{E}_x[x | y]] \quad (\text{As } \text{var}_y[y] = \mathbf{E}_y[y^2] - (\mathbf{E}_y[y])^2)
\end{aligned}$$