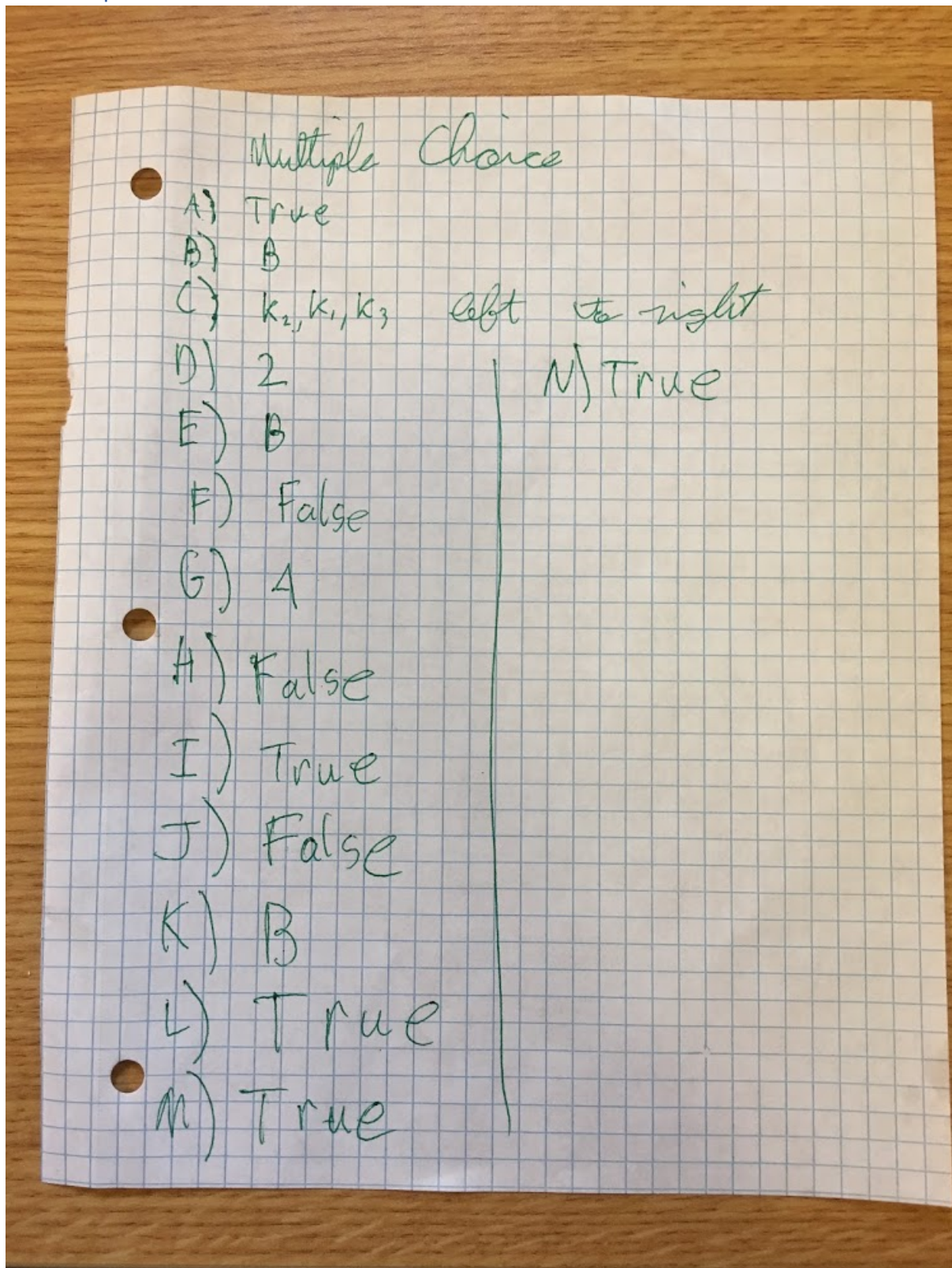


1 Multiple Choice:



2 Naïve Bayes

Q1

	P
Happy = YES	$\frac{1}{2}$
Happy = NO	$\frac{1}{2}$

	yr = FRESH	Grade = 4
Happy → yes	$\frac{1}{3}$	$\frac{2}{3}$
Happy → No	$\frac{1}{2}$	$\frac{1}{3}$

Q2: apply chain rule

$$P(\text{yr} = \text{Fresh}, \text{Grade} = C, \text{Happy} = \text{NO})$$

$$= P(\text{happy} = \text{NO}) \cdot P(\text{Grade} = C \mid \text{happy} = \text{NO})$$

$$\begin{aligned} & \cdot P(\text{yr} = \text{Fresh} \mid \text{happy} = \text{NO}) \quad \frac{1+2}{2+4} = \frac{1}{2} \\ & = \frac{1+3}{2+4} \cdot \frac{1}{2} \cdot \frac{2}{3} \cdot \frac{1}{2} = \frac{1}{6} \end{aligned}$$

Part 3

```
Var h = random()
  if h > P(Happy=Yes)
    Happy = NO
  else
    Happy = Yes
← Var y = random()
  if y > P(yr = Fresh | Happy)
    Year = Senior
  else
    Year = Fresh
Var g = random
  if g > P(A | HAPPY)
    Grade = C
  else
    Grade = A
```

Return ~~(Happy)~~
y(Grade, Year, Happy)

~~return (H~~

3 Data Transformation

Part 1/2

① $w^T x = \tilde{w}^T \tilde{x}$ 3 Data Trans

$$w^T x = \tilde{w}^T (Ax)$$

$$\cancel{w^T x} \xrightarrow{1} \tilde{w}^T A = \tilde{w} A^T$$

$$\cancel{w^T x} \xrightarrow{2} \tilde{w}^T A$$

$$\cancel{w^T} = \cancel{A^T} \tilde{w}$$

transpose

$$w = \tilde{w} A^T$$

$$w = \tilde{w} A$$

$$\Rightarrow \tilde{w} = A^{-1} w$$

②

$$\argmin_w \left(\argmin_{\tilde{w}} \frac{\lambda}{2} \|A^{-1} w\|^2 + \sum_i (y_i - w^T x_i)^2 \right)$$

$$\argmin_w \frac{\lambda}{2} \|A^{-1} w\|^2 + \sum_i (y_i - w^T x_i)^2$$

Part 3

We know that our answer from part b

$$\begin{aligned} & \arg \min_w \frac{\lambda}{2} \|A^{-1}w\|^2 + \sum_i (y_i - w^T x_i)^2 \\ &= A \left[\arg \min_{\tilde{w}} \frac{\lambda}{2} \|\tilde{w}\|^2 + \sum_i (y_i - \tilde{w}^T \tilde{x}_i)^2 \right] \end{aligned}$$

What is effectively going on is that the data is transformed into a new rescaled space and then ridge regression is performed and then the data is transformed back. Since the transformation in question is simply a rescaling, all that transforming really means is multiplying by a diagonal matrix A . Because of this, the model parameter, w can either be optimized directly by training the first objective or \tilde{w} can be optimized and then transformed after the fact. The point is basically that transforming in this case is a well behaved 1 to 1 linear scaling transform within the same space that doesn't really change the result of the optimizing problem. It is possible that such a scaling could help avoid numerical issues in practice, however, apart from that, this transformed model class is exactly equally expressive as regular ridge regression and the model class isn't changed in any functional way.

4 Latent Markov Embedding

Question 1

We begin by noting that the dual point model is no less expressive than the single point model.

We can see this simply by noting that dual point is a generalization of single point or single point is a special case of dual point (the thinking is equivalent).

Consider some matrix X as in a single point model optimum. We can see that a dual point model with $U = V = X$ is the same model as the single point model and therefore the dual point model as a whole can't have a $P(S)$ less than that of the single point model.

Worst case the two are just the same, otherwise dual point will train to a greater maximum. Thus, we have shown $P(S)$ for the dual-point model is *never less* than $P(S)$ for the single point model.

Question 2

This would imply that we are dealing with the case in which single point and dual point result in the same model

$$\Rightarrow U = V = X$$

5 Neural Net

Parts 1 and 2

5 Neural Net

1) $\frac{\partial}{\partial w_{11}} (y - f(x))^2 = \frac{\partial (y - f(x))^2}{\partial f} \cdot \frac{\partial f}{\partial w_{11}}$

$f(x) = \sigma(u_1 h_1(x) + u_2 h_2(x))$

$\sigma(w_{11}x_1 + w_{21}x_2)$ $\sigma(w_{12}x_1 + w_{22}x_2)$

$f(x) = \sigma(u_1 \underbrace{\sigma(w_{11}x_1 + w_{21}x_2)}_{I_1^*} + u_2 \underbrace{\sigma(w_{12}x_1 + w_{22}x_2)}_{I_2^*})$

$\frac{\partial f}{\partial w_{11}} = \sigma'(I) \cdot I' = \frac{\partial \sigma(I)}{\partial I} \cdot \frac{\partial I}{\partial w_{11}}$ (constant WRT w_{11})

$= \sigma(I)(1 - \sigma(I)) \cdot \frac{\partial I}{\partial w_{11}}$

$\frac{\partial I}{\partial w_{11}} = \frac{\partial I_1^*}{\partial w_{11}} = \frac{\partial \sigma(I_1^*)}{\partial I_1^*} \cdot \frac{\partial I_1^*}{\partial w_{11}}$

$= u_1 \sigma(I_1^*)(1 - \sigma(I_1^*)) \cdot x_1$

$\frac{\partial f}{\partial w_{11}} = \sigma(I)(1 - \sigma(I)) \cdot [\sigma(I_1^*)(1 - \sigma(I_1^*)) \cdot x_1] u_1$

$$\frac{\partial}{\partial w_{11}} (y - f(x))^2 = -2(y - f(x)) \cdot \sigma(I) (1 - \sigma(I))$$

$$\cdot \left[\sigma(I^*) (1 - \sigma(I^*)) u_1 x_1 \right] u_1$$

for I and I^* as defined above
for given values

$$u_1 \sigma(I_1^*) = 0.5 \sigma(0.25 \cdot 0.1 + 0.05 \cdot 0.5)$$

$$= 0.2565$$

$$u_2 \sigma(I_2^*) = -0.1 \sigma(0.1 \cdot 0.1 + (-0.25) \cdot 0.5)$$

$$= -0.047128$$

$$\Rightarrow I = 0.2565 + -0.047128 = 0.20937$$

$$\Rightarrow f(x) = f(I) = 0.552153 = \sigma(I)$$

$$\frac{\partial L}{\partial w_{11}} = -2(0.75 - 0.552153) \cdot (0.552153)(1 - 0.552153) \cdot$$

$$(0.512497)(1 - 0.512497) 0.1 \cdot 0.5 = \boxed{-0.00122}$$

Part 3

We have:

$$\frac{\partial}{\partial w_{11}} L(y, f(x)) = \frac{\partial}{\partial w_{11}} (y - f(x))^2 = \frac{\partial}{\partial f} (y - f(x))^2 \cdot \frac{\partial f}{\partial w_{11}}$$

We can clearly see that the term $\frac{\partial f}{\partial w_{11}}$ is the culprit behind vanishing gradients. Calculating this gradient involves calculating, and multiplying by through chain rule, gradients of composite functions from all previous layers, so we have a lot of multiplied terms that can cause the result gradient to vanish. Additionally, we are using

$$\sigma(s) = \frac{e^s}{1 + e^s}$$

as our non-linear activation function

and we have that $\max_s \sigma'(s) = \frac{1}{4} < 1$, so in the process of calculating $\frac{\partial f}{\partial w_{11}}$, we must also calculate and multiply by many such small values which is exacerbated by having more layers and thus more derivatives of activation functions.