The honor code.

(a). Zheng Wu

this hw is quite time consuming. Dealing with dummy variables takes a large amount of time.

(b). I hereby state that all of my solutions were entirely in my words and were written by me. I have not looked at other student's solutions and I have fairly credited all external sources in this write up.

hw6 sec154.

1. Computational complexity.
   ① a+v    $O(d)$
       $a^T v$    $O(d)$

   ② A+B.   $O(n \cdot d)$.
      space storing matrix A.   $O(n \cdot d)$

   ③. A v.    $O(n \cdot d^2)$.
      $A^T B$    $O(d \cdot n^2)$.

   ④. $A^T B v$

2. kernel methods.
   ① computing $X^T X$ takes ~~off(t)~~ $O(l^2 n)$ .
      inverse matrix takes. $O(l^3)$
        thus adding up these leading terms.
     computing $\hat{\theta} x$ takes $O(l^3 + n \cdot l^2)$ complexity.

   ②. Using raw data matrix X.
   original form $w^* = (X^T X + \lambda I)^{-1} X^T y$
       $(X^T X + \lambda I) w^* = X^T y$
       $X^T X w^* + \lambda w^* = X^T y$
         $\lambda w^* = X^T y - X^T X \cdot w^*$.
         $w^* = \dfrac{X^T (y - X \cdot w^*)}{\lambda}$    which is the linear combination of columns of $X^T$.
     i.e. $W = X^T \cdot v$

     substitute into   $X^T X w + \lambda w = X^T y$
        $\Rightarrow$   $X^T X (X^T v) + \lambda (X^T v) = X^T y$
          $X^T (X X^T v + \lambda v) = X^T y$

   we cannot cancel $X^T$ directly. but if $X X^T v + \lambda v = y$ has solutions

it implies that $X^T(XX^Tv+\lambda v)=X^Ty$ also has solutions.

$$XX^Tv+\lambda v=y$$
$$v^* = (XX^T+\lambda)^{-1}y.$$

substitute $v^*$ back into $w=X^Tv$

$$\Rightarrow w= X^T(XX^T+\lambda)^{-1}y$$

i.e. $\hat{\theta}_\lambda = X^T(XX^T+\lambda I_n)^{-1}y$ is an equivalent & valued estimate

for ridge problem.

complexity:    computing $XX^T$ takes $O(n^2l)$

              inverting matrix takes $O(n^3)$

              thus adding up these two leading terms. we get $O(n^3+n^2l)$

③   For $\hat{\theta}_\lambda = (X^TX+\lambda I)^{-1}X^Ty \cdots O(l^3+l^2n)$

     for $\hat{\theta}_\lambda = X^T(XX^T+\lambda I_n)^{-1}y \cdots O(n^3+n^2l)$

④.   we extend the feature map.

A feature map is used when the problem cannot be
linearly solved. feature map $\phi$, map data to a new
space where a linear model could be applied to solve
the original learning problem.

In addition, kernel function is defined to map $\phi$ to a higher
dimensional space and compute inner-product based similarity
$$\phi(x)^T\phi(z) \text{ in that space}.$$

⑤

(a).