

(1)

- (A) True \rightarrow some λ exists where L1-regression regularized has same solution as given constrained Regression
- (B) D $\rightarrow (1 - \frac{1}{n}) \rightarrow P(\text{No } x, y) \rightarrow N$ samples per bootstrap, M bootstraps $A^{M \cdot N}$
- (C) k_1, k_3, k_2
- (D) $\hat{y} = 4$
- (E) A \Rightarrow same application of weight w in (2) encourages correlation
- (F) False, \Rightarrow Hard margin SVM just won't converge not $w=0$
- (G) A \Rightarrow support vectors remain the same, so no change in w, b
- (H) False \Rightarrow new predictions can change to incorrect or not change previous predictions
- (I) True \Rightarrow linear algebra operations provide closed form solution with Z matrices fixed
- (J) False \Rightarrow Bias variance tradeoff shows test error not minimized with var or bias min.
- (K) B $\Rightarrow \prod_{j=1}^M P(x^j | y^j) P(y^j | y^{j-2})$ same for all values, each model learns equal
- (L) True \Rightarrow exists a tradeoff between accuracy and size/variability of training space
- (M) True \Rightarrow With infinite data any classifier function can be learned with decision trees
- (N) True \Rightarrow Infinite data means ~~and~~ infinitely harder for model to fit on specific points and not ~~generalize~~ generalize

NAIVE BAYES

$$(1) P(\text{Grade} = A | \text{Happy} = \text{Yes}) = 2/3 \quad P(\text{Grade} = A | \text{Happy} = \text{No}) = 2/3$$

$$P(\text{Grade} = C | \text{Happy} = \text{Yes}) = 2/3 \quad P(\text{Grade} = C | \text{Happy} = \text{No}) = 2/3$$

$$P(\text{Year} = \text{Freshman} | \text{Happy} = \text{Yes}) = 2/3 \quad P(\text{Year} = \text{Freshman} | \text{Happy} = \text{No}) = 1/2$$

$$P(\text{Year} = \text{Senior} | \text{Happy} = \text{Yes}) = 2/3 \quad P(\text{Year} = \text{Senior} | \text{Happy} = \text{No}) = 1/2$$

$$P(\text{Happy} = \text{Yes}) = 1/2 \quad P(\text{Happy} = \text{No}) = 1/2$$

$$(2) P(y, g, H) = P(H) \cdot P(y|H) \cdot P(g|H)$$

$$= P(\text{Happy} = \text{No}) \cdot P(\text{Freshman} | \text{Not happy}) \cdot P(C | \text{Not happy})$$

$$= 1/2 \cdot 1/2 \cdot 2/3 = 2/12 = 1/6$$

$$P(\text{Year} = \text{Freshman}, \text{Grade} = C, \text{Happy?} = \text{No}) = 2/6$$

$$(3) \text{ if random() } > P(\text{Happy?} = \text{No})$$

Random(.) represents generating probabilities via random number [0, 1]

~~happy = True~~

~~happy = False~~

happy = True

True = Yes, False = No

else:

happy = False

$$\text{ if random() } > P(\text{Grade} = C | \text{Happy?} = \text{happy})$$

grade = A

Probabilities can be called from training results

else:

grade = C

$$\text{ if random() } > P(\text{Year} = \text{Freshman} | \text{Happy?} = \text{happy})$$

year ~~grade~~ = Senior

else:

year ~~grade~~ = Freshman

return happy, grade, ~~fresh~~ year

(3) DATA TRANSFORMATIONS

(1) $w^T x = \tilde{w}^T \tilde{x}$, given $\tilde{x} = Ax$ we can substitute

$$w^T x = \tilde{w}^T Ax \text{ which implies } w^T = \tilde{w}^T A$$

$$\Rightarrow (w^T)^T = (\tilde{w}^T A)^T \Rightarrow w = A^T \tilde{w}$$

So we get w as a function of A and \tilde{w} where:

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

$$(2) \text{ given } w = A^T \tilde{w} \Rightarrow \tilde{w} = (A^T)^{-1} w \Rightarrow \tilde{w}^T = w^T A^{-1}$$

$$\Rightarrow \arg \min_{\tilde{w}} \frac{\lambda}{2} \|\tilde{w}\|^2 + \sum_i (y_i - \tilde{w}^T \tilde{x}_i)^2$$

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

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

$$\text{optimization problem: } \arg \min_w \frac{\lambda}{2} \|(A^T)^{-1} w\|^2 + \sum_i (y_i - w^T x_i)^2$$

(3) Notice from (2) the ^{squared loss} ~~expression~~ ^{squared} term simplifies to the

standard ridge regression loss term. Thus, the only difference

is in the regularization term where $(A^T)^{-1} w$ is L2 regularized

instead of just w . In context of A as a transformation on x ,

this essentially means the optimization from (2) changes from

(1) in that the ~~se~~ inverted scaling of the weights is regularized

compared to the standard ridge.

(4)

(1) Notice that U and V can always converge to x where the dual point model now has $u, v = x$ in its optimization. Since this optimization is always a choice we can argue that the dual point model in its worst case, can always do as well as the best performance of the single-point model since the dual point can always replicate the single point. As a result, given optimal u, v, x we ~~can~~ can conclude the dual-point ~~model~~ model likelihood can only be as worse as single-point model, so the dual-point $p(s)$ likelihood is never less than the single point model.

(2) If the transition probabilities of both single-point (8) and dual point (6) models are the same it implies that $U=V$ since $u(s') - v(s') = x(s') - x(s)$ for all s .

(5)

$$(A) \text{ Chain rule: } \frac{\partial}{\partial w_{11}} L(y, f(x)) = \frac{\partial (y - f(x))^2}{\partial w_{11}}$$

$$= \frac{\partial (y - f(x))^2}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial (\sum u_i h_i(x))} \cdot \frac{\partial (\sum u_i h_i(x))}{\partial h_1(x)} \cdot \frac{\partial h_1(x)}{\partial \sum w_{ji} x_j} \cdot \frac{\partial \sum w_{ji} x_j}{\partial w_{11}}$$

$$= \left[\underbrace{-2(y - f(x))}_{\sigma(s) \text{ derivative}} \right] \underbrace{\left[f(x)(1 - f(x)) \right]}_{\sigma(s) \text{ derivative}} \left[u_1 \right] \left[h_1(x)(1 - h_1(x)) \right] \left[x_1 \right]$$

$$(B) \text{ } h_1(x) = \sigma(w_{11}x_1 + w_{21}x_2) = \sigma(.05) = .5124$$

$$h_2(x) = \sigma(w_{12}x_1 + w_{22}x_2) = \sigma(-.115) = .4713$$

$$f(x) = \sigma(u_1 h_1(x) + u_2 h_2(x)) = \sigma((.5124)(.5) - (.4713)(.2)) = \sigma(.209) = .5522$$

$$\Rightarrow -2(.75 - .5522)(.5522(1 - .5522))(.5)(.5124(1 - .5124))(.2)$$

$$\Rightarrow -.002222 \Rightarrow \left[\frac{\partial}{\partial w_{11}} L(y, f(x)) = -.002222 \right]_{\text{done}}$$

(C) The sigmoid term as well as its derivatives (which include sigmoid terms) result in the vanishing gradient problem, notice

in (A) there are $f(x)$ and $h_1(x)$ terms which are both

non-linearities, non-linearities always result in output < 1 which

when multiplied ⁱⁿ to chain rule will go to zero resulting in a

vanishing gradient. For larger layer networks this chain rule is larger

resulting in more multiplications and likelihood of vanishing, exacerbating the problem