

• 1. a) $FPR = 0.2 = \frac{FP}{N=30}$ 50 misclassified
 ~~$FNR = 1 - 0.2 = 0.8$~~ $\Rightarrow P = 200 - 30 = 170$

$$FP = 30(0.2) = \frac{30}{5} = 6 \Rightarrow TN = 30 - 6 = 24$$

$$\Rightarrow FN = 50 - 6 = 44$$

$$\Rightarrow TP = 170 - 44 = 126$$

$$\boxed{TPR = \frac{126}{170}}$$

b) $FNR = 1 - TPR$

$$= 1 - \frac{126}{170}$$

$$= \boxed{\frac{44}{170}}$$

c) $acc = \frac{TP + TN}{P + N} = \frac{126 + 24}{200} = \frac{150}{200} = \boxed{\frac{3}{4}}$

d) $F1 = \frac{2 \cdot \text{prec} \cdot \text{recall}}{\text{prec} + \text{recall}}$

$$= \boxed{\frac{2 \left(\frac{126}{132} \cdot \frac{126}{170} \right)}{\frac{126}{132} + \frac{126}{170}}}$$

$$\text{prec} = \frac{TP}{P} = \frac{126}{TP + FP}$$

$$= \frac{126}{126 + 6} = \frac{126}{132}$$

$$\text{recall} = TPR = \frac{126}{170}$$

2.

$$X = \begin{pmatrix} 2 & 3 \\ -1 & 4 \\ 3 & 0 \\ 2 & -1 \end{pmatrix}$$

$$X^T = \begin{pmatrix} 2 & -1 & 3 & 2 \\ 3 & 4 & 0 & -1 \end{pmatrix}$$

Gram

$$XX^T = \begin{pmatrix} 13 & 10 & 6 & 1 \\ 10 & 17 & -3 & -6 \\ 6 & -3 & 9 & 6 \\ 1 & -6 & 6 & 5 \end{pmatrix}$$

3. Training data

X is 100×5
row column

where each row is a data point
with 5 features

y is 100×1 corresponding to
the labels of each data point

w is a 5×1 (if we make
the data zero mean)
vector representing
the decision boundary

ϵ is the residual or error
in our solution
 100×1 vector

4. Sensitivity to outliers means our decision boundary will ~~screw~~ towards outliers, and likely not find a "perfect" boundary for linearly separable data.

we can either try to penalize points with large margins

use SVM / soft margin SVM (in case of non separable data)

or use perceptron (potentially with a limit on the number of iterations if the data is not linearly separable)

5. we may want the cluster exemplars to be data points from the training data. The method is the point with the minimal avg dissimilarity and is by design guaranteed to be a training point

6. soft margin SVM works when the data is not linearly separable. It introduces a slack variable for each data point allowing for margin errors

7. margin

$$z(x) = \frac{y(w^T x - t)}{\|w\|}$$

$$w = (3, 2, -1)$$

$$t = 4$$

$$\|w\| = \sqrt{9+4+1} = \sqrt{14}$$

a) $x = (0, 0, 0)$ $y_i = -1$

$$z(x) = \frac{-1(0-4)}{\sqrt{14}} = \boxed{\frac{4}{\sqrt{14}}}$$

b) $x = (2, 1, 1)$ $y_i = -1$

$$z(x) = \frac{-1(6+2-1-4)}{\sqrt{14}} = \frac{-3}{\sqrt{14}} \quad \text{misclassified}$$

c) $x = (1, 2, 3)$ $y_i = 1$

$$z(x) = \frac{(3+4-3-4)}{\sqrt{14}} = 0 \quad \text{on the boundary}$$

8. minimize within-class variance
maximize between-class variance

9. a) $\boxed{H4}$

we choose hypotheses that maximize (ML)
the likelihood ($P(D|H)$)

b) choose H that maximizes a
posteriori ($P(H|D)$)

$$P(H|D) = \frac{P(H) P(D|H)}{P(D) = 0.37}$$

this is $\boxed{H3}$

10. Chebyshev distance

$$L_{\infty}(x, y) = \|x - y\|_{\infty} = \max_i |x_i - y_i|$$

points

$(1, 2), (1, 1), (2, 1), (3, 2), (4, 2), (4, 3)$

dist
 $(2, 3) : \downarrow \quad \downarrow$
 $1, 2, 2, 1, 2, 2$

$(4, 1.5) : 3, 3, 2, 1, 0.5, 1.5$

iter 0

cluster 1: mean: $(2, 3)$: data pts: $(1, 2), (1, 1), (2, 1), (3, 2)$
 cluster 2: pts: $(4, 2), (4, 3)$

iter 1

cluster 1 mean

$$\frac{(6, 6)}{4} = (1.5, 1.5) \text{ pts. } (1, 2), (1, 1), (2, 1), (3, 2)$$

cluster 2 mean

$$\frac{(8, 5)}{2} = (4, 2.5) : \text{pts } (4, 2), (4, 3)$$

dist:

$(1.5, 1.5) \quad 0.5 \quad 0.5 \quad 0.5 \quad 0.5 \quad 1.5 \quad 1.5$

$(4, 2.5) \quad 3 \quad 3 \quad 2 \quad 1 \quad \underbrace{0.5 \quad 0.5}_{(4, 2) \quad (4, 3)}$

iter 2

same as iter 1

11. we propagate error backwards in order to update weight connecting nodes, ultimately these weights are what gets output.

12. 3 hidden layers

$$4 \times 3 + 3 \times 3 + 3 \times 2 = 12 + 9 + 6 = \boxed{27} \text{ weights}$$

$$\begin{array}{cccc} 0 & & & \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array}$$

13. run multiple times with different initial weights and pick the best one
smallest error on validation set

14. $(1.5, 2.0, 1.0) \cdot (0.5, -0.3, 0.2)$

$$0.75 + (-0.6) + 0.2 = 0.25$$

squared error

$$(0.5 - 0.25)^2 = \boxed{(0.25)^2}$$

15. Intrinsic dimensionality is 4

16. to transfer data into a higher dimensional space that is linearly separable

17. d could appear $\left\lfloor \frac{|D|}{N} \right\rfloor$ times

18. 2 errors / 5 = 0.4

$\epsilon = 0.4$ (I realize this isn't quite right)

new weights

$$\left\{ \underbrace{\frac{0.2}{2\epsilon} \mid \frac{0.1}{2\epsilon}}_{\text{increase}} \mid \underbrace{\frac{0.35}{2(1-\epsilon)} \mid \frac{0.3}{2(1-\epsilon)} \mid \frac{0.05}{2(1-\epsilon)}}_{\text{decrease}} \right\}$$

$$1000/10 = 100$$

19. each training set has
900 data pts

any given data points d_i will appear in
the test set exactly once

20. convolution layer has

$32 \times 11 \times 11 \times 3$ weights to learn
pooling layer

$$2 \times 2 \times 2$$

+ 54×54 activation maps

+ $64 \times 64 \times 3$ input size