

Final Pset

Saturday, March 23, 2019 8:38 PM

1. The polynomial transform of order $Q = 10$ applied to \mathcal{X} of dimension $d=2$ results in a \mathcal{Z} space of what dimensionality (not counting the constant term $x_0 = 1$ or $z_0 = 1$)?
- [a] 12
 - [b] 20
 - [c] 35**
 - [d] 100
 - [e] None of the above**

55. I thought we were approximating...

The feature space X has two dimensions. For a 10th order polynomial, not including constants, we have the following terms:

We can express the "activation" of a z vector as follows

$$\sum_{i=1}^{10} \sum_{j=0}^i a_{ij} x_1^i x_2^{j-i}$$
$$\sum_{i=1}^{10} (i+1) = 10 + \sum_{i=1}^{10} i = 10 + 11 * \frac{10}{2} = 65$$

$$10+9+8+7+6+5+4+3+2+1 = 55$$

2. Recall that the average hypothesis \bar{g} was based on training the model g on different data sets \mathcal{D} to get $g^{(\mathcal{D})} \in \mathcal{H}$, and taking the average w.r.t. \mathcal{D} to get \bar{g} . Which of the following models \mathcal{H} does \bar{g} belong to?

~~[a] A singleton \mathcal{H} (\mathcal{H} has one hypothesis)~~

~~[b] \mathcal{H} is the set of all constant, real-valued hypotheses~~

~~[c] \mathcal{H} is the linear regression model~~

[d] \mathcal{H} is the logistic regression model

~~[e] None of the above~~

ion $d = 2$ re-
nt coordinate

ms:

training the same model \mathcal{H}
the expected value of $g^{(\mathcal{D})}$
could result in $\bar{g} \notin \mathcal{H}$?

ses

[a] It is the logistic regression model

- [e] None of the above

All hypothesis spaces are "convex spaces"

The "average" of any set constants is a constant.

The average of a set of linear models can be arranged into another linear model where the weight vector is the average (commutative and associative properties repeatedly)

For the log-reg models, this is not the case.

• Overfitting

3. Which of the following statements is *false*?

- [a] If there is overfitting, there must be two or more hypotheses that differ in their different values of E_{in} . True! Overfitting is part of learning, which is discriminating hypotheses.
- [b] If there is overfitting, there must be two or more hypotheses that differ in their different values of E_{out} . True! Again, we are learning in a way that degrades E_{out} .
- [c] If there is overfitting, there must be two or more hypotheses that differ in their different values of $(E_{out} - E_{in})$. $E_{out2} > E_{out1}, E_{in2} < E_{in1}, E_{out2} - E_{in2} > E_{out1} - E_{in1}$
we can re-arrange to put c_1 and c_2 on the same side of the equation
- [d] We can always determine if there is overfitting by comparing the values of $(E_{out} - E_{in})$. Sure, but if we had that so readily available.....
- [e] We cannot determine overfitting based on one hypothesis only.
true

Overfitting in general, is the specific phenomenon of "learning" in a way that we decrease E_{in} , but are increasing E_{out} .
problem/function we are trying to learn about.

4. Which of the following statements is true?

- [a] Deterministic noise cannot occur with stochastic noise
- [b] Deterministic noise does not depend on the hypothesis set
- [c] Deterministic noise does not depend on the target function
- [d] Stochastic noise does not depend on the hypothesis set.
- [e] Stochastic noise does not depend on the target distribution

age weight vector of the linear model (leverage

have

theses.

have

have

$E_{out2} - E_{in2} = E_{out1} + c_1 - E_{in1} + c_2$. c_1 and c_2 are non-negative, therefore

and c_2 on one side

es of

We do this by "learning the dataset" and not the

et. Deterministic noise is another word for bias. Bias is the inability of a Hset to represent the target function.

on.



ion. Don't know what this means.

● Regularization

5. The regularized weight \mathbf{w}_{reg} is a solution to:

$$\text{minimize } \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \text{ subject to } \mathbf{w}^T \Gamma^T$$

where Γ is a matrix. If $\mathbf{w}_{\text{lin}}^T \Gamma^T \Gamma \mathbf{w}_{\text{lin}} \leq C$, where \mathbf{w}_{lin} is the solution, then what is \mathbf{w}_{reg} ?

[a] $\mathbf{w}_{\text{reg}} = \mathbf{w}_{\text{lin}}$

If \mathbf{w}_{lin} is vanilla lin reg solution, and it meets our regularization constraint

[b] $\mathbf{w}_{\text{reg}} = \Gamma \mathbf{w}_{\text{lin}}$

Then we done! Our "simple" / "reduced" hypothesis set is consistent

[c] $\mathbf{w}_{\text{reg}} = \Gamma^T \Gamma \mathbf{w}_{\text{lin}}$

Turns out, our concerns over stochastic noise were great overestimated

[d] $\mathbf{w}_{\text{reg}} = C \Gamma \mathbf{w}_{\text{lin}}$

[e] $\mathbf{w}_{\text{reg}} = C \mathbf{w}_{\text{lin}}$

6. Soft-order constraints that regularize polynomial models can be

~~[a] written as hard-order constraints~~

~~[b] translated into augmented error~~ 

~~[c] determined from the value of the VC dimension~~ 

~~[d] used to decrease both E_{in} and E_{out}~~ More constraint, E_{in} is non-decreasing!

[e] None of the above is true

● Bayesian Priors

19. Let $f \in [0, 1]$ be the unknown probability of getting a heart attack in a certain population. Notice that f is just a constant, not a function of time or other variables.

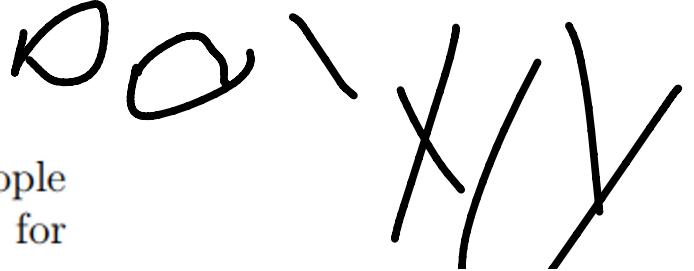
$${}^T \Gamma \mathbf{w} \leq C,$$

the linear regression

ization constraint....

still GOT THE JOB DONE.

ly exaggerated :)

> 
for people
unction, for

19. Let $J \in [0, 1]$ be the unknown probability of getting a near attack in a certain population. Notice that f is just a constant, not a function of simplicity. We want to model f using a hypothesis $h \in [0, 1]$. Before we have any data, we assume that $P(h = f)$ is uniform over $h \in [0, 1]$ (the prior). Suppose we sample one person from the population, and it turns out that he or she has had a near attack. Which of the following is true about the posterior probability to f given this sample point?

- [a] The posterior is uniform over $[0, 1]$.
- [b] The posterior increases linearly over $[0, 1]$.
- [c] The posterior increases nonlinearly over $[0, 1]$.
- [d] The posterior is a delta function at 1 (implying f has to be 1).
- [e] The posterior cannot be evaluated based on the given information.**

Probability of 0 is 0

• Aggregation

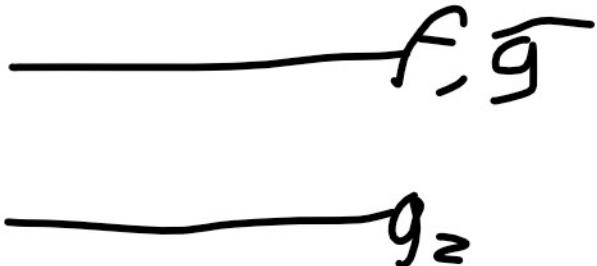
20. Given two learned hypotheses g_1 and g_2 , we construct the aggregate hypothesis g given by $g(\mathbf{x}) = \frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x}))$ for all $\mathbf{x} \in \mathcal{X}$. If we use mean-squared error, which of the following statements is true?

- [a] $E_{\text{out}}(g)$ cannot be worse than $E_{\text{out}}(f)$. Why not?
- [b] $E_{\text{out}}(g)$ cannot be worse than the average of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$.
- [c] $E_{\text{out}}(g)$ cannot be worse than the average of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$.**
- [d] $E_{\text{out}}(g)$ has to be between $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$ (including the end values of that interval).
- [e] None of the above

Before we had a flat thing.

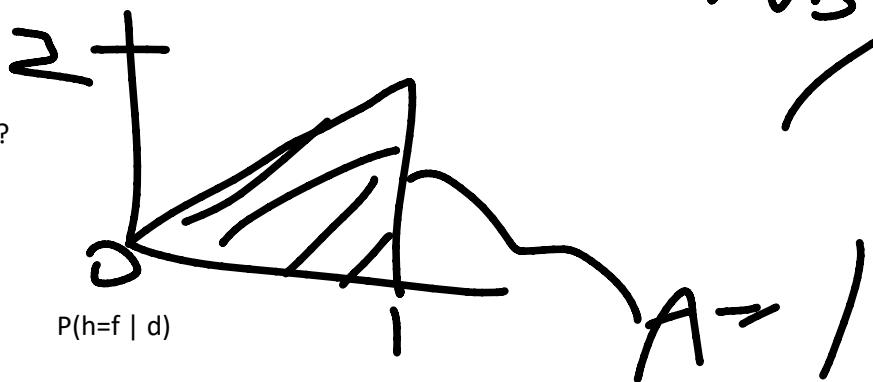
Now, we learn on $|D| = 1$. Specifically, we know $P(h=f|D) = P(h=f)$ is uniform and $P(D|f)$ is constant throughout, we can reformat the integral and simply compute then "normalize" via integration. The shape of the potential "constant value" of f is linear asap.)

by nature
of SS



for people
ction, for
we see any
We pick
ad a heart
that $h = f$

1 2
controversy



ion.

$P(D|h=f)*P(h=f) / P(D)$. Given our prior is
so $P(h=f|D) = \alpha P(D|h=f)$, which, we
of this "posterior distribution" of the

\bar{g}
 \bar{g}_2

7. Set $\lambda = 1$ and do not apply a feature transform (i.e., use $\mathbf{z} = \mathbf{x} = (1, x_1, x_2)$). Which among the following classifiers has the lowest E_{in} ?
- [a] 5 versus all
 - [b] 6 versus all
 - [c] 7 versus all
 - [d] 8 versus all**
 - [e] 9 versus all
8. Now, apply a feature transform $\mathbf{z} = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$, and set $\lambda = 1$. Which among the following classifiers has the lowest E_{out} ?
- [a] 0 versus all
 - [b] 1 versus all**
 - [c] 2 versus all
 - [d] 3 versus all
 - [e] 4 versus all
9. If we compare using the transform versus not using it, and apply that to '0 versus all' through '9 versus all', which of the following statements is correct for $\lambda = 1$?
- [a] Overfitting always occurs when we use the transform. False, some improve
 - [b] The transform always improves the out-of-sample performance by at least 5% (E_{out} with transform $\leq 0.95E_{\text{out}}$ without transform). False, some impr
 - [c] The transform does not make any difference in the out-of-sample performance.

Definitively untrue

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{z}_n - y_n)^2 + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

$$\begin{aligned}
E_{\text{in}}(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \\
&= \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \\
&= \frac{1}{N} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}),
\end{aligned}$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}). \quad +2\lambda/N \mathbf{w}^T = 0$$

above some dont

(3.3)

(3.4)

- [d] The transform always worsens the out-of-sample performance by at least 5%.
- [e] The transform improves the out-of-sample performance of ‘5 versus all,’ but by less than 5%.
10. Train the ‘1 versus 5’ classifier with $\mathbf{z} = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$ with $\lambda = 0.01$ and $\lambda = 1$. Which of the following statements is correct?
- [a] Overfitting occurs (from $\lambda = 1$ to $\lambda = 0.01$).
- [b] The two classifiers have the same E_{in} .
- [c] The two classifiers have the same E_{out} .
- [d] When λ goes up, both E_{in} and E_{out} go up.
- [e] When λ goes up, both E_{in} and E_{out} go down.

● Support Vector Machines

11. Consider the following training set generated from a target function $f : \mathcal{X} \rightarrow \{-1, +1\}$ where $\mathcal{X} = \mathbb{R}^2$

$$\begin{array}{lll} \mathbf{x}_1 = (1, 0), y_1 = -1 & \mathbf{x}_2 = (0, 1), y_2 = -1 & \mathbf{x}_3 = (0, -1), y_3 = -1 \\ \mathbf{x}_4 = (-1, 0), y_4 = +1 & \mathbf{x}_5 = (0, 2), y_5 = +1 & \mathbf{x}_6 = (0, -2), y_6 = +1 \\ \mathbf{x}_7 = (-2, 0), y_7 = +1 & & \end{array}$$

Transform this training set into another two-dimensional space \mathcal{Z}

$$z_1 = x_2^2 - 2x_1 - 1 \quad z_2 = x_1^2 - 2x_2 + 1$$

Using geometry (not quadratic programming), what values of \mathbf{w} (without w_0) and b specify the separating plane $\mathbf{w}^T \mathbf{z} + b = 0$ that maximizes the margin in the \mathcal{Z} space? The values of w_1, w_2, b are:

- [a] $-1, 1, -0.5$
- [b] $1, -1, -0.5$
- [c] $1, 0, -0.5$
- [d] $0, 1, -0.5$
- [e] None of the above would work.

12. Consider the same training set of the previous problem but instead of explicitly

12. Consider the same training set of the previous problem, but instead of explicitly transforming the input space \mathcal{X} , apply the hard-margin SVM algorithm with the kernel

$$K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$

(which corresponds to a second-order polynomial transformation). Set up the expression for $\mathcal{L}(\alpha_1 \dots \alpha_7)$ and solve for the optimal $\alpha_1, \dots, \alpha_7$ (numerically, using a quadratic programming package). The number of support vectors you get is in what range?

- [a] 0-1
- [b] 2-3
- [c] 4-5
- [d] 6-7
- [e] >7

● Radial Basis Functions

We experiment with the RBF model, both in regular form (Lloyd + pseudo-inverse) with K centers:

$$\text{sign} \left(\sum_{k=1}^K w_k \exp(-\gamma \|\mathbf{x} - \mu_k\|^2) + b \right)$$

(notice that there is a bias term), and in kernel form (using the RBF kernel in hard-margin SVM):

$$\text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n \exp(-\gamma \|\mathbf{x} - \mathbf{x}_n\|^2) + b \right).$$

The input space is $\mathcal{X} = [-1, 1] \times [-1, 1]$ with uniform probability distribution, and the target is

$$f(\mathbf{x}) = \text{sign}(x_2 - x_1 + 0.25 \sin(\pi x_1))$$

which is slightly nonlinear in the \mathcal{X} space. In each run, generate 100 training points at random using this target, and apply both forms of RBF to these training points. Here are some guidelines:

- Repeat the experiment for as many runs as needed to get the answer to be stable (statistically away from flipping to the closest competing answer).
- In case a data set is not separable in the ' \mathcal{Z} space' by the RBF kernel using hard-margin SVM, discard the run but keep track of how often this happens, if ever.

-
- When you use Lloyd's algorithm, initialize the centers to random points in \mathcal{X} and iterate until there is no change from iteration to iteration. If a cluster becomes empty, **discard the run and repeat**.

- 13.** For $\gamma = 1.5$, how often do you get a data set that is not separable by the RBF kernel (using hard-margin SVM)? *Hint: Run the hard-margin SVM, then check that the solution has $E_{\text{in}} = 0$.*

- [a] $\leq 5\%$ of the time
 [b] $> 5\%$ but $\leq 10\%$ of the time
 [c] $> 10\%$ but $\leq 20\%$ of the time
 [d] $> 20\%$ but $\leq 40\%$ of the time
 [e] $> 40\%$ of the time

\mathbb{F}_k

14. If we use $K = 9$ for regular RBF and take $\gamma = 1.5$, how often does the kernel form beat the regular form (excluding runs mentioned in Problem 13 and runs with empty clusters, if any) in terms of E_{out} ?

- [a] $\leq 15\%$ of the time
- [b] $> 15\%$ but $\leq 30\%$ of the time
- [c] $> 30\%$ but $\leq 50\%$ of the time
- [d] $> 50\%$ but $\leq 75\%$ of the time
- [e] $> 75\%$ of the time**

X
15. If we use $K = 12$ for regular RBF and take $\gamma = 1.5$, how often does the kernel form beat the regular form (excluding runs mentioned in Problem 13 and runs with empty clusters, if any) in terms of E_{out} ?

- [a] $\leq 10\%$ of the time
- [b] $> 10\%$ but $\leq 30\%$ of the time
- [c] $> 30\%$ but $\leq 60\%$ of the time
- [d] $> 60\%$ but $\leq 90\%$ of the time**
- [e] $> 90\%$ of the time**

l
s

el
s

16. Now we focus on regular RBF only, with $\gamma = 1.5$. If we go from $K = 9$ clusters to $K = 12$ clusters (only 9 and 12), which of the following 5 cases happens most often in your runs (excluding runs with empty clusters, if any)? Up or down means strictly so.

- [a] E_{in} goes down, but E_{out} goes up.

7

- [b] E_{in} goes up, but E_{out} goes down.
- [c] Both E_{in} and E_{out} go up.
- [d] Both E_{in} and E_{out} go down.
- [e] E_{in} and E_{out} remain the same.

17. For regular RBF with $K = 9$, if we go from $\gamma = 1.5$ to $\gamma = 2$ (only 1.5 and 2), which of the following 5 cases happens most often in your runs (excluding runs with empty clusters, if any)? Up or down means strictly so.

- [a] E_{in} goes down, but E_{out} goes up.
- [b] E_{in} goes up, but E_{out} goes down.
- [c] Both E_{in} and E_{out} go up.
- [d] Both E_{in} and E_{out} go down.
- [e] E_{in} and E_{out} remain the same.

s
t
n

),
as

18. What is the percentage of time that regular RBF achieves $E_{\text{in}} = 0$ with $K = 9$ and $\gamma = 1.5$ (excluding runs with empty clusters, if any)?

- [a] $\leq 10\%$ of the time
- [b] $> 10\%$ but $\leq 20\%$ of the time
- [c] $> 20\%$ but $\leq 30\%$ of the time
- [d] $> 30\%$ but $\leq 50\%$ of the time
- [e] $> 50\%$ of the time

Hard-Margin SVM with Kernel

1: Construct Q_D from the kernel K , and A_D :

$$Q_D = \begin{bmatrix} y_1 y_1 K_{11} & \dots & y_1 y_N K_{1N} \\ y_2 y_1 K_{21} & \dots & y_2 y_N K_{2N} \\ \vdots & \vdots & \vdots \\ y_N y_1 K_{N1} & \dots & y_N y_N K_{NN} \end{bmatrix} \quad \text{and} \quad A_D = \begin{bmatrix} \mathbf{y}^T \\ -\mathbf{y}^T \\ \mathbf{I}_{N \times N} \end{bmatrix},$$

where $K_{mn} = K(\mathbf{x}_m, \mathbf{x}_n)$. (K is called the *Gram* matrix.)

2: Use a QP-solver to optimize the dual problem:

$$\alpha^* \leftarrow \text{QP}(Q_D, -\mathbf{1}_N, A_D, \mathbf{0}_{N+2}).$$

3: Let s be any support vector for which $\alpha_s^* > 0$. Compute

$$b^* = y_s - \sum_{\alpha_n^* > 0} y_n \alpha_n^* K(\mathbf{x}_n, \mathbf{x}_s).$$

4: Return the final hypothesis

$$g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n^* > 0} y_n \alpha_n^* K(\mathbf{x}_n, \mathbf{x}) + b^* \right).$$

x

