

1. For a regression problem with 2 features, consider the effect of different regularizers and different amounts of regularization, graphically as described below. You may do this by hand, or you may use a computer to assist you if you prefer.

Assume the unconstrained objective function is $f_{obj}(\underline{w}) = \frac{1}{N} \text{RSS}(\underline{w}, \mathcal{D}_i)$. For simplicity, in this problem we assume $w_0 = 0$, consistent with a dataset that has been standardized in both x and y . Consider 10 different datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{10}$, each resulting in an unconstrained (unregularized) minimum at $\hat{\underline{w}}_{\text{lin}}^{(i)}$ given by:

$$\begin{aligned} \mathcal{D}_1 : \hat{\underline{w}}_{\text{lin}}^{(1)} &= (10, 0), & \mathcal{D}_2 : \hat{\underline{w}}_{\text{lin}}^{(2)} &= (10, 2), & \mathcal{D}_3 : \hat{\underline{w}}_{\text{lin}}^{(3)} &= (10, 4), \\ \mathcal{D}_4 : \hat{\underline{w}}_{\text{lin}}^{(4)} &= (10, 6), & \mathcal{D}_5 : \hat{\underline{w}}_{\text{lin}}^{(5)} &= (8, 6), & \mathcal{D}_6 : \hat{\underline{w}}_{\text{lin}}^{(6)} &= (8, 8), \\ \mathcal{D}_7 : \hat{\underline{w}}_{\text{lin}}^{(7)} &= (6, 8), & \mathcal{D}_8 : \hat{\underline{w}}_{\text{lin}}^{(8)} &= (6, 10), & \mathcal{D}_9 : \hat{\underline{w}}_{\text{lin}}^{(9)} &= (4, 10), & \mathcal{D}_{10} : \hat{\underline{w}}_{\text{lin}}^{(10)} &= (2, 10) \end{aligned}$$

Assume the shape of $\text{RSS}(\underline{w}, \mathcal{D}_i) = \text{constant}$ curves in 2D weight space are circles (special case of ellipses), for simplicity.

In each regularizer case given below, make a plot in 2D weight space, showing:

- (i) the 10 unregularized-minimum points $\hat{\underline{w}}_{\text{lin}}^{(i)}$ given above,
- (ii) the region that satisfies the given regularizer constraint, and
- (iii) the resulting 10 regularized minimum points, i.e., solution of

$$\hat{\underline{w}}_{\text{reg}}^{(i)} = \arg \min_{\underline{w}} f_{obj}(\underline{w}, \mathcal{D}_i) \quad \text{s.t.} \quad \Omega(\underline{w}) \leq C.$$

for each i . Also show or justify how you found the resulting $\hat{\underline{w}}_{\text{reg}}^{(i)}$. (Showing your method for one or two points in each regularizer case, should be sufficient.)

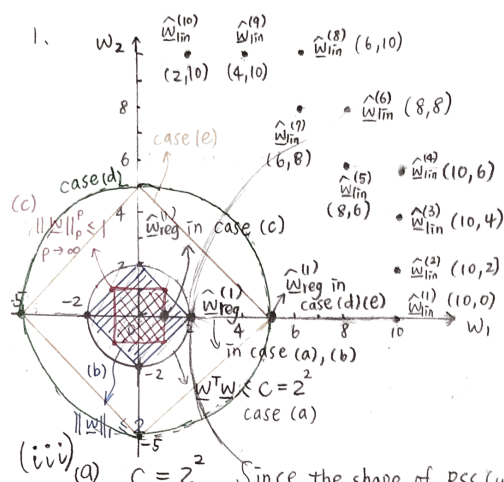
- (iv) Also, answer: how many of the resulting $\hat{\underline{w}}_{\text{reg}}^{(i)}$, $i = 1, 2, \dots, 10$, are more sparse than the corresponding $\hat{\underline{w}}_{\text{lin}}^{(i)}$? For the purpose of this problem, define sparsity as the number of components $\hat{w}_j^{(i)}$ that have value 0, for a given i .

Tip: For cases in which there are more than one possible $\hat{\underline{w}}_{\text{reg}}^{(i)}$ for a given dataset and a given constraint, pick any one.

- (a) L2 regularization: $\Omega(\underline{w}) = \|\underline{w}\|_2^2$, $C = 2^2$.
- (b) L1 regularization: $\Omega(\underline{w}) = \|\underline{w}\|_1$, $C = 2$.
- (c) L_p regularization (based on p -norm): $\Omega(\underline{w}) = \|\underline{w}\|_p^p$, as $p \rightarrow \infty$, $C = 1$.

Hint: if you're not sure of the shape of $\|\underline{w}\|_p^p = 1$, try plotting it numerically for increasing p , e.g. $p = 4, 10, 100$.

- (d) Repeat (a), except with $C = 5^2$.
- (e) Repeat (b), except with $C = 5$.



$$\begin{aligned}
 f_{obj}(w) &= \frac{1}{N} \text{RSS}(w, \mathcal{D}_i) = \frac{1}{N} \sum_{i=1}^N (y - w^T x_i)^2 \\
 &= \frac{1}{N} (y - Xw)^T (y - Xw) \\
 &= \frac{1}{N} (y^T y - 2y^T Xw + w^T X^T X w) \\
 \Rightarrow \nabla_w f_{obj}(w) &= \frac{1}{N} (2X^T y - 2X^T X w) = 0 \\
 \Rightarrow \hat{w}_{lin} &= (X^T X)^{-1} X^T y = \bar{X} \bar{y} \\
 \Rightarrow \bar{X} \bar{y} &= (\bar{X} \bar{X}) \hat{w}_{lin}
 \end{aligned}$$

(iii) (a) $C = 2^2$. Since the shape of $\text{RSS}(w, \mathcal{D}_i) = \text{constant}$ curves are circles, $\Omega(w) = \|w\|_2^2$ the regularized minimum point will on the line passing through original point and $\hat{w}_{lin}^{(i)}$. Also, the point will on the circle $w^T w = 2^2$.

① find $\hat{w}_{reg}^{(1)}$

$\hat{w}_{reg}^{(1)}$ is on the line $y=0$ and the circle $w^T w = 2^2 \Rightarrow \hat{w}_{reg}^{(1)} = (2, 0)$

② find $\hat{w}_{reg}^{(2)}$

$$y = ax + b \begin{cases} 10a + b = 2 \\ b = 0 \end{cases} \Rightarrow y = \frac{1}{5}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{1}{25}x^2 = 4 \Rightarrow x = 5\sqrt{\frac{2}{13}} \Rightarrow \hat{w}_{reg}^{(2)} = (5\sqrt{\frac{2}{13}}, \sqrt{\frac{2}{13}})$$

③ find $\hat{w}_{reg}^{(3)}$

$$y = ax + b \begin{cases} 10a + b = 4 \\ b = 0 \end{cases} \Rightarrow y = \frac{2}{5}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{4}{25}x^2 = 4 \Rightarrow x = 10\sqrt{\frac{1}{29}} \Rightarrow \hat{w}_{reg}^{(3)} = (\frac{10}{\sqrt{29}}, \frac{4}{\sqrt{29}})$$

④ find $\hat{w}_{reg}^{(4)}$

$$\Rightarrow y = \frac{3}{5}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{9}{25}x^2 = 4 \Rightarrow x = 5\sqrt{\frac{2}{17}} \Rightarrow \hat{w}_{reg}^{(4)} = (5\sqrt{\frac{2}{17}}, 3\sqrt{\frac{2}{17}})$$

⑤ find $\hat{w}_{reg}^{(5)}$

$$\Rightarrow y = \frac{3}{4}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{9}{16}x^2 = 4 \Rightarrow x = \frac{8}{5} \Rightarrow \hat{w}_{reg}^{(5)} = (\frac{8}{5}, \frac{6}{5})$$

⑥ find $\hat{w}_{reg}^{(6)}$

$$\Rightarrow y = x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + x^2 = 4 \Rightarrow x = \sqrt{2} \Rightarrow \hat{w}_{reg}^{(6)} = (\sqrt{2}, \sqrt{2})$$

⑦ find $\hat{w}_{reg}^{(7)}$

$$\Rightarrow y = \frac{4}{3}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{16}{9}x^2 = 4 \Rightarrow x = \frac{6}{5} \Rightarrow \hat{w}_{reg}^{(7)} = (\frac{6}{5}, \frac{8}{5})$$

⑧ find $\hat{w}_{reg}^{(8)}$

$$\Rightarrow y = \frac{5}{3}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{25}{9}x^2 = 4 \Rightarrow x = 3\sqrt{\frac{2}{17}} \Rightarrow \hat{w}_{reg}^{(8)} = (3\sqrt{\frac{2}{17}}, 5\sqrt{\frac{2}{17}})$$

⑨ find $\hat{w}_{reg}^{(9)}$

$$\Rightarrow y = \frac{5}{2}x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + \frac{25}{4}x^2 = 4 \Rightarrow x = 4\sqrt{\frac{1}{29}} \Rightarrow \hat{w}_{reg}^{(9)} = (\frac{4}{\sqrt{29}}, \frac{10}{\sqrt{29}})$$

⑩ find $\hat{w}_{reg}^{(10)}$

$$\Rightarrow y = 5x$$

$$x^2 + y^2 = 4 \Rightarrow x^2 + 25x^2 = 4 \Rightarrow x = \sqrt{\frac{2}{13}} \Rightarrow \hat{w}_{reg}^{(10)} = (\sqrt{\frac{2}{13}}, 5\sqrt{\frac{2}{13}})$$

(b) $C=2$, $\Omega(\underline{w}) = \|\underline{w}\|$,

Only when the point $\hat{\underline{w}}_{lin}^{(i)}$ is located on the region between the line $y=x+2$ and $y=x-2$, the $\hat{\underline{w}}_{reg}^{(i)}$ will locate on the line $y=-x+2$ which tangent to the circle $RSS(\underline{w}, \mathcal{D}_i) = \text{constant}$. If not, the $\hat{\underline{w}}_{reg}^{(i)}$ will locate on $(2,0)$ or $(0,2)$

$$\hat{\underline{w}}_{reg}^{(1)} = (2,0), \hat{\underline{w}}_{reg}^{(2)} = (2,0), \hat{\underline{w}}_{reg}^{(3)} = (2,0), \hat{\underline{w}}_{reg}^{(4)} = (2,0), \hat{\underline{w}}_{reg}^{(5)} = (2,0)$$

$$\hat{\underline{w}}_{reg}^{(7)} = (0,2), \hat{\underline{w}}_{reg}^{(8)} = (0,2), \hat{\underline{w}}_{reg}^{(9)} = (0,2), \hat{\underline{w}}_{reg}^{(10)} = (0,2)$$

$$\hat{\underline{w}}_{reg}^{(6)} = (1,1)$$

(c) $C=1$, $\Omega(\underline{w}) = \|\underline{w}\|_p^p$ as $p \rightarrow \infty$

If $\hat{\underline{w}}_{lin}^{(i)}$ is located in the region between $y=1$ and $y=-1$, $\hat{\underline{w}}_{reg}^{(i)}$ will be on the line $x=1$

If $\hat{\underline{w}}_{lin}^{(i)}$ is located in the region between $x=-1$ and $x=1$, $\hat{\underline{w}}_{reg}^{(i)}$ will be on the line $y=1$.

If $\hat{\underline{w}}_{lin}^{(i)}$ is located outside of the region mentioned above, $\hat{\underline{w}}_{reg}^{(i)} = (1,1)$

$$\hat{\underline{w}}_{reg}^{(1)} = (1,0), \hat{\underline{w}}_{reg}^{(2)} = (1,1), \hat{\underline{w}}_{reg}^{(3)} = (1,1), \hat{\underline{w}}_{reg}^{(4)} = (1,1), \hat{\underline{w}}_{reg}^{(5)} = (1,1)$$

$$\hat{\underline{w}}_{reg}^{(6)} = (1,1), \hat{\underline{w}}_{reg}^{(7)} = (1,1), \hat{\underline{w}}_{reg}^{(8)} = (1,1), \hat{\underline{w}}_{reg}^{(9)} = (1,1), \hat{\underline{w}}_{reg}^{(10)} = (1,1)$$

(d) $C=5^2$, $\Omega(\underline{w}) = \|\underline{w}\|_2^2$

① $\hat{\underline{w}}_{reg}^{(1)} = (5,0)$

② $y = \frac{1}{5}x$

$$x^2 + \frac{1}{25}x^2 = 25 \Rightarrow x = \frac{25}{\sqrt{26}} \Rightarrow \hat{\underline{w}}_{reg}^{(2)} = \left(\frac{25}{\sqrt{26}}, \frac{5}{\sqrt{26}} \right)$$

③ $y = \frac{2}{5}x$

$$x^2 + \frac{4}{25}x^2 = 25 \Rightarrow x = 25/\sqrt{29} \Rightarrow \hat{\underline{w}}_{reg}^{(3)} = \left(\frac{25}{\sqrt{29}}, \frac{10}{\sqrt{29}} \right)$$

④ $y = \frac{3}{5}x$

$$x^2 + \frac{9}{25}x^2 = 25 \Rightarrow x = 25/\sqrt{34} \Rightarrow \hat{\underline{w}}_{reg}^{(4)} = \left(\frac{25}{\sqrt{34}}, \frac{15}{\sqrt{34}} \right)$$

⑤ $y = \frac{4}{5}x$

$$x^2 + \frac{16}{25}x^2 = 25 \Rightarrow x = 4 \Rightarrow \hat{\underline{w}}_{reg}^{(5)} = (4, 3)$$

⑥ $y = x$

$$x^2 + y^2 = 25 \Rightarrow x^2 + x^2 = 25 \Rightarrow x = \frac{5}{\sqrt{2}} \Rightarrow \hat{\underline{w}}_{reg}^{(6)} = \left(\frac{5}{\sqrt{2}}, \frac{5}{\sqrt{2}} \right)$$

⑦ $y = \frac{4}{3}x$

$$x^2 + \frac{16}{9}x^2 = 25 \Rightarrow x = 3 \Rightarrow \hat{\underline{w}}_{reg}^{(7)} = (3, 4)$$

⑧ $y = \frac{5}{3}x$

$$x^2 + \frac{25}{9}x^2 = 25 \Rightarrow x = \frac{15}{\sqrt{34}} \Rightarrow \hat{\underline{w}}_{reg}^{(8)} = \left(\frac{15}{\sqrt{34}}, \frac{25}{\sqrt{34}} \right)$$

⑨ $y = \frac{5}{2}x$

$$x^2 + \frac{25}{4}x^2 = 25 \Rightarrow x = \frac{10}{\sqrt{29}} \Rightarrow \hat{\underline{w}}_{reg}^{(9)} = \left(\frac{10}{\sqrt{29}}, \frac{25}{\sqrt{29}} \right)$$

⑩ $y = 5x$

$$x^2 + 25x^2 = 25 \Rightarrow x = \frac{5}{\sqrt{26}} \Rightarrow \hat{\underline{w}}_{reg}^{(10)} = \left(\frac{5}{\sqrt{26}}, \frac{25}{\sqrt{26}} \right)$$

(e) $c=5, \Omega(w) = \|w\|_1$

Only when the point $\hat{w}_{lin}^{(i)}$ is located on the region between the line $y=x+5$ and $y=x-5$, the $\hat{w}_{reg}^{(i)}$ will locate on the line $y=-x+5$ which tangent to the circle

$RSS(w, \mathcal{D}_i) = \text{constant}$. If not, the $\hat{w}_{reg}^{(i)}$ will locate on $(5,0)$ or $(0,5)$

$$\hat{w}_{reg}^{(1)} = (5,0), \hat{w}_{reg}^{(2)} = (5,0), \hat{w}_{reg}^{(3)} = (5,0), \hat{w}_{reg}^{(9)} = (0,5), \hat{w}_{reg}^{(10)} = (0,5).$$

$$\hat{w}_{reg}^{(4)}: (x, -x+5) = (\frac{9}{2}, \frac{1}{2})$$

$$\left(\begin{aligned} (x-5, -x+5) \cdot (10-x, 6-(-x+5)) &= (x-\frac{5}{2})(10-x) + (-x+\frac{5}{2})(x+1) = 0 \\ \Rightarrow x &= \frac{9}{2} \text{ or } \cancel{x}, y = \frac{1}{2} \text{ or } \emptyset \end{aligned} \right)$$

$$\hat{w}_{reg}^{(5)}: (x, -x+5) = (\frac{7}{2}, \frac{3}{2})$$

$$\left(\begin{aligned} (x-5, -x+5) \cdot (8-x, x+1) &= (x-5)(8-x) + (-x+5)(x+1) = 0 \\ \Rightarrow x &= \frac{7}{2} \text{ or } \cancel{x}, y = \frac{3}{2}, \emptyset \end{aligned} \right)$$

$$\hat{w}_{reg}^{(6)}: (x, -x+5) = (\frac{5}{2}, \frac{5}{2})$$

$$\left(\begin{aligned} (x-5, -x+5) \cdot (8-x, x+3) &= 0 \Rightarrow (x-5)(8-x) + (-x+5)(x+3) = 0 \\ \Rightarrow x &= \frac{5}{2} \text{ or } \cancel{x}, y = \frac{5}{2} \text{ or } \emptyset \end{aligned} \right)$$

$$\hat{w}_{reg}^{(7)}: (x, -x+5) = (\frac{3}{2}, \frac{7}{2})$$

$$\left(\begin{aligned} (x-5, -x+5) \cdot (6-x, x+3) &= 0 \Rightarrow (x-5)(6-x) + (-x+5)(x+3) = 0 \\ \Rightarrow x &= \frac{3}{2} \text{ or } \cancel{x}, y = \frac{7}{2} \text{ or } \emptyset \end{aligned} \right)$$

$$\hat{w}_{reg}^{(8)}: (x, -x+5) = (\frac{1}{2}, \frac{9}{2})$$

$$\left(\begin{aligned} (x-5, -x+5) \cdot (6-x, x+5) &= 0 \Rightarrow (x-5)(6-x) + (-x+5)(x+5) = 0 \\ \Rightarrow x &= \frac{1}{2} \text{ or } \cancel{x}, y = \frac{9}{2} \text{ or } \emptyset \end{aligned} \right)$$

(iv) (a) In this case, no resulting $\hat{w}_{reg}^{(i)}$ are more sparse than the corresponding $\hat{w}_{lin}^{(i)}$

(b) In this case, $\hat{w}_{reg}^{(i)}, i=2, 3, 4, 5, 7, 8, 9, 10$ (8 $\hat{w}_{reg}^{(i)}$ in total) become more sparse than $\hat{w}_{lin}^{(i)}$.

(c) In this case, no $\hat{w}_{reg}^{(i)}$ becomes more sparse

(d) In this case, no $\hat{w}_{reg}^{(i)}$ becomes more sparse

(e) In this case, $\hat{w}_{reg}^{(i)}, i=2, 3, 9, 10$ (4 $\hat{w}_{reg}^{(i)}$ in total) become more sparse than $\hat{w}_{lin}^{(i)}$.

2. Suppose you develop and optimize a machine learning system, starting with setting aside a test dataset \mathcal{D}_{Test} , and using the remaining data points as the set \mathcal{D}' . Your hypothesis set is \mathcal{H}_1 , and you use \mathcal{D}' as a training set to find its best hypothesis h_{g1} . Let $d_{VC}(\mathcal{H}_1) = d_{VC}^{(1)}$, $N' = |\mathcal{D}'|$, and $N_{Test} = |\mathcal{D}_{Test}|$. When you are finished, you pull out the test set and calculate $E_{Test}(h_{g1})$. In this problem, all generalization bounds are with tolerance δ (with probability $\geq 1 - \delta$).
- Draw a flow chart (like we did in Lecture 16, p. 6, and like AML Fig. 4.11), that shows the dataset usage, hypothesis set, and procedure.
 - Give an inequality for the generalization bound based on the training error $E_{\mathcal{D}'}(h_{g1})$, and the generalization bound based on the test-set error $E_{Test}(h_{g1})$.

Afterwards, independently of the results you got above, you think of a different approach that you also want to try. So you start the process all over again, setting aside the same test set \mathcal{D}_{Test} . You define a hypothesis set \mathcal{H}_2 for your model. Let $d_{VC}(\mathcal{H}_2) = d_{VC}^{(2)}$.

In this case, however, you also use some model selection to choose the optimum number of features in a feature selection process. So you split \mathcal{D}' into a training set \mathcal{D}_{Tr} and a validation set \mathcal{D}_{Val} , that are disjoint. You use \mathcal{D}_{Tr} to train each model (based on a given number of features d), and use model selection to compare different values of d , with $d = 1, 2, 3, \dots, d_{\max}$, in which d_{\max} is the maximum number of features you try. You choose the best number of features by comparing $E_{Val}(h_{g2}^{(d)})$ for each value of d . Let $N_{Tr} = |\mathcal{D}_{Tr}|$, and $N_{Val} = |\mathcal{D}_{Val}|$.

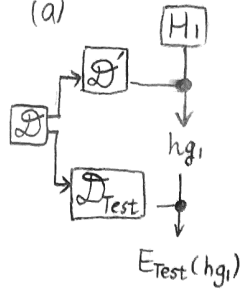
- Draw a flow chart (like we did in Lecture 16, p. 6, and like AML Fig. 4.11), that shows the dataset usage, hypothesis sets, parameter values d , and procedure, for this second approach only.
- Give:
 - An inequality for the generalization bound based on the training-set error $E_{Tr}(h_{g2}^{(d)})$ for a given number of features d ;
 - An inequality for the generalization bound based on the validation-set error $E_{Val}(h_{g2}^{(d^*)})$ for the optimal number of features d^* ;
 - An inequality for the generalization bound based on the test-set error $E_{\mathcal{D}_{Test}}(h_{g2}^{(d^*)})$ for the best hypothesis $h_{g2}^{(d^*)}$.

Finally, you compare the best results from the 2 systems you developed, and pick the one with the lower test-set error.

- Give an inequality for the generalization bound based on the test-set error $E_{Test}(h_g^*)$ for the best hypothesis h_g^* .

Hint: what is the effective hypothesis set used by \mathcal{D}_{Test} to pick between the two machine-learning systems you developed?

2. (a)

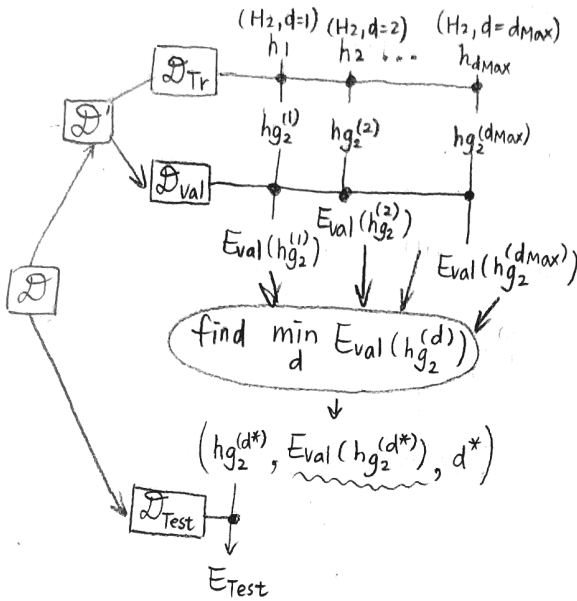


$$(b) E_{out}(hg_1) \leq E_{D'}(hg_1) + \sqrt{\frac{8}{N} \ln \frac{4[(2N)^{d_{vc}^2} + 1]}{\delta}}$$

$$E_{out}(hg_1) \leq E_{Test}(hg_1) + \sqrt{\frac{1}{2N_{Test}} \ln \frac{2}{\delta}}$$

both with probability $\geq 1 - \delta$

(c)



(d)

$$(i) E_{out}(hg_2^{(d)}) \leq E_{Tr}(hg_2^{(d)}) + \sqrt{\frac{8}{N_{Tr}} \ln \frac{4[(2N_{Tr})^{d_{vc}^{(2)}} + 1]}{\delta}}$$

$$(ii) E_{out}(hg_2^{(d*)}) \leq E_{val}(hg_2^{(d*)}) + \sqrt{\frac{1}{2N_{val}} \ln \frac{2 \cdot d_{max}}{\delta}}$$

$$(iii) E_{out}(hg_2^{(d*)}) \leq E_{Test}(hg_2^{(d*)}) + \sqrt{\frac{1}{2N_{Test}} \ln \frac{2}{\delta}}$$

All of (i), (ii), (iii) with probability $\geq 1 - \delta$

(e)

When we compare the best results from the 2 systems we developed, it means $\mathcal{H} = \{hg_1, hg_2^{(d*)}\} \Rightarrow$ we have 2 hypotheses

$\Rightarrow M = 2$

$$E_{out}(hg^*) \leq E_{Test}(hg^*) + \sqrt{\frac{1}{2N_{Test}} \ln \frac{2 \times 2}{\delta}} \text{ with probability } \geq 1 - \delta$$