

CS 156a Final

Claire Goeckner-Wald

December 2, 2016

Nonlinear Transforms

1. [e] None of the above.

The polynomial transform of order $Q = 10$ applied to χ of dimension $d = 2$ results in a \mathbb{Z} space of dimensionality $\tilde{d} = 65$, not counting the constant coordinates $x_0 = 1$ and $z_0 = 1$. The polynomial transform of order Q applied to some $x \in \chi = (x_1, x_2)$ consists of all variations of $x_1^i x_2^j$ with non-negative integers i and j satisfying $i + j \leq Q$. For example, the polynomial transform of order $Q = 3$ applied to χ of dimension $d = 2$ is

$$\begin{aligned} z &= (x_1^0 x_2^0, x_1^1 x_2^0, x_1^0 x_2^1, x_1^1 x_2^1, x_1^2 x_2^0, x_1^0 x_2^2, x_1^3 x_2^0, x_1^0 x_2^3); \\ &= (1, x_1, x_2, x_1 x_2, x_1^2 x_2, x_1^2 x_2, x_1^3, x_2^3). \end{aligned}$$

In this case, the dimensionality, \tilde{d} , of the \mathbb{Z} -space is 7. (We do not count the constant coordinate z_0 .) Then, for some polynomial transform of order q , we have $\tilde{d} = (\sum_{i=0}^q (q - i) + 1) - 1$. Therefore, for $Q = 10$, we have $\tilde{d} = (\sum_{i=0}^{10} (10 - i) + 1) - 1$, or $\tilde{d} = 65$.

Bias and Variance

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

[a] If \mathcal{H} is a singleton, then the chosen hypothesis $g^{\mathcal{D}}$ will be the same, no matter the data set \mathcal{D} . Thus, when finding the average value of $g^{\mathcal{D}}$ with respect to \mathcal{D} , the average hypothesis must be the single hypothesis in \mathcal{H} . Therefore, $\bar{g} \in \mathcal{H}$.

[b] If \mathcal{H} is the set of all constant, real-valued hypotheses, then the chosen hypothesis $g^{\mathcal{D}}$ for some data \mathcal{D} will be a constant real number $b \in \mathbb{R}$. Therefore, when finding the average value of $g^{\mathcal{D}}$ with respect to \mathcal{D} , the expected value of \bar{g} must also be some real-valued constant $\bar{b} \in \mathbb{R}$. Thus, since $\bar{b} \in \mathbb{R}$, then \mathcal{H} must contain \bar{b} . Therefore, $\bar{g} \in \mathcal{H}$.

[c] If \mathcal{H} is the linear regression model, then it represents hypotheses consisting of a polynomial (or, binomial) with real-valued coefficients. Therefore, when finding the average value of $g^{\mathcal{D}}$ with respect to \mathcal{D} , the expected value of \bar{g} must also be some set of real-valued coefficients. Thus, since that set is real-valued, then \mathcal{H} must contain it as a hypothesis. Therefore, $\bar{g} \in \mathcal{H}$.

[d] If \mathcal{H} is the logistic regression model, then it represents hypotheses as sigmoidal functions. Then, when finding the chosen hypothesis $g^{\mathcal{D}}$ for some data \mathcal{D} , the hypothesis is a sigmoid ranging from 0 to 1. However, when finding the average value of $g^{\mathcal{D}}$ with respect to \mathcal{D} , the expected value of \bar{g} is not necessarily a sigmoid. Given sigmoids $f(x) = \frac{e^{a+bx}}{1+e^{a+bx}}$ and $f(x) = \frac{e^{a+bx}}{1+e^{a+bx}}$, then $h(x) = f(x) + g(x)$ is not necessarily a sigmoid. Since \mathcal{H} contains sigmoids, \bar{g} is not necessarily in \mathcal{H} .

Overfitting

3. [d] is false

[a] *True*: Overfitting is the result of picking a hypothesis with a higher E_{out} because of a misleadingly low E_{in} . Therefore, there must be two or more hypotheses with differing values of E_{in} , in order for us to be misled into overfitting.

[b] *True*: If there is overfitting, there must be two or more hypotheses that have different values of E_{out} because overfitting results in choosing an incorrect hypothesis g with a higher E_{out} as a cause of a misleadingly lower E_{in} . Thus, there must be two or more hypotheses with different values of E_{out} in order for overfitting to occur.

[c] *True*: If there is overfitting, there must be two or more hypotheses that have different values of $(E_{out} - E_{in})$. For example, consider the set \mathcal{H} composed of hypotheses h_1 and h_2 , each with out-sample errors $e_{out,1}$ and $e_{out,2}$ respectively, where $e_{out,1} < e_{out,2}$. Therefore, it is best if we choose that $g = h_1$. Overfitting occurs when h_2 , the inferior hypothesis, has a misleading low value of E_{in} . Let each hypothesis have in-sample errors $e_{in,1}$ and $e_{in,2}$. Then, in the case of overfitting, $e_{in,1} > e_{in,2}$. Thus, in the case of overfitting, we have $(e_{out,1} - e_{in,1}) < (e_{out,2} - e_{in,2})$. Therefore, if there is overfitting then two or more hypotheses that have different values of $(E_{out} - E_{in})$.

[d] *False*: We cannot always determine if there is overfitting using values of $(E_{out} - E_{in})$. Consider hypothesis set \mathcal{H} with two hypotheses, h_1 and h_2 . In this example, we have no overfitting because we select hypothesis h_1 for its superior E_{in}

$$e_{in,1} = 0.2, e_{out,1} = 0.5, (e_{out,1} - e_{in,1}) = 0.3e_{in,2} = 0.6, e_{out,2} = 0.7, (e_{out,2} - e_{in,2}) = 0.1$$

Now, consider hypothesis set \mathcal{H}' with two hypotheses, h'_1 and h'_2 . In this example, assuming a naive learning algorithm, we have overfitting because we are misled by the misleading in-sample error of h'_1 .

$$e_{in,1} = 0.2, e_{out,1} = 0.5, (e_{out,1} - e_{in,1}) = 0.3e_{in,2} = 0.3, e_{out,2} = 0.4, (e_{out,2} - e_{in,2}) = 0.1$$

Since we are only given $(E_{out} - E_{in})$, we will similar values for both scenarios- however, one demonstrates overfitting, and the other does not. Thus, we cannot always determine the presence of overfitting given $(E_{out} - E_{in})$.

[e] *True*: We cannot determine overfitting based on one hypothesis only.

4. [d] is true

[a] *False*: Deterministic noise can occur with stochastic noise- the two are not existentially dependent.

[b] *False*: Deterministic noise depends on the hypothesis set, because we can construct certain hypothesis sets \mathcal{H} that are prone to high deterministic noise. For example, if a hypothesis set \mathcal{H}_2 , composed of exclusively second-order polynomials, attempts to fit a tenth-order polynomial target function with noise, then the deterministic noise will likely be high.

[c] *False*: Deterministic noise depends on the target function, as seen in the lower right-hand graph of the logo of this class. In this graph, the increasing target complexity results in a worse E_{out} , due to deterministic noise. For example, if we have a hypothesis set \mathcal{H}_2 , composed of exclusively second-order polynomials, attempting to fit a variety of different n -order target functions, some target functions will result in a lower deterministic noise than others.

[d] *True*: Stochastic noise does not depend on the hypothesis set. Stochastic noise is a byproduct of the random fluctuations of the given data. For example, stochastic noise may occur where human error is present in measurements performed to create the training data set.

[e] *False*: Stochastic noise does depend on the target function. Stochastic noise is random fluctuations of given data set that is formed using the target function. Thus, if the target function changes, we may have different stochastic noise.

Regularization

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

In the case that the solution for \mathbf{w}_{lin} is already constrained by $\mathbf{w}^\top \Gamma^\top \Gamma \mathbf{w} \leq C$, then we simply use that value of \mathbf{w}_{lin} as our regularized weight, \mathbf{w}_{reg} . This is because \mathbf{w}_{lin} is also found by minimizing $\frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$, so if \mathbf{w}_{lin} satisfies the constraint as well, then $\mathbf{w}_{\text{reg}} = \mathbf{w}_{\text{lin}}$.

6. [b] translated into augmented error

[a] Writing a soft-order constraint as a hard-order constraint is not the point of the soft-order constraint.

[b] Soft-order constraints that regularize polynomial models can be translated into augmented error through a series of calculations.

$$\text{Minimizing } E_{\text{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^\top (\mathbf{Z}\mathbf{w} - \mathbf{y}) \text{ subject to } \mathbf{w}^\top \mathbf{w} \leq C$$

Can be solved by $\mathbf{w}_{\text{reg}} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{y}$, with λ being the “amount of regularization”

[c] The value of the VC *dimension* does not determine the soft-order constraint. However, the VC *formulation* is equivalent to a soft-order constraint.

[d] Using soft-order constraints may result in an *increased* E_{in} and decreased E_{out} simultaneously. (This is kind of the goal of soft-order constraints.)

Regularized Linear Regression

- 7.
- 8.
- 9.
- 10.

Support Vector Machines

- 11.
- 12.

Radial Basis Functions

- 13.
- 14.
- 15.
- 16.
- 17.
- 18.

Bayesian Priors

- 19.

Aggregation

- 20. [c] $E_{\text{out}}(g)$ cannot be worse than the average of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$ *is true* **NOTE: Expand on [c]**

[a] *False:* $E_{\text{out}}(g)$ can be worse than $E_{\text{out}}(g_1)$. Consider the case where g_2 is consistently farther from the target $y = f(\mathbf{x})$ for $\mathbf{x} \in \chi$ than g_1 is. Then, taking the average of outputs, $E_{\text{out}}(g)$, would be greater than $E_{\text{out}}(g_1)$.

[b] *False*: $E_{\text{out}}(g)$ can be worse than the minimum of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$. Consider the minimum to be that of g_{\min} , which is either g_1 or g_2 . The other shall be g_{\max} . In this case, g_{\min} , will, by definition, have a lower $E_{\text{out}}(g_{\min})$ than that of g_{\max} . Thus, g_{\max} tends to be farther from the target $y = f(\mathbf{x})$ for $\mathbf{x} \in \chi$ than g_{\min} is. Thus, $E_{\text{out}}(g)$ would be greater than $E_{\text{out}}(g_{\min})$.

[c] *True*: $E_{\text{out}}(g)$ cannot be worse than the average of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$ because we are using mean squared error. This means that $E_{\text{out}}(g)$, even if higher than $E_{\text{out}}(g_1)$ or $E_{\text{out}}(g_2)$, cannot be worse than their average.

$$\begin{aligned}
E_{\text{out}}(g) &= \mathbb{E}[(f(\mathbf{x}) - g(\mathbf{x}))^2] = \mathbb{E}\left[\left(f(\mathbf{x}) - \frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x}))\right)^2\right] \\
\frac{1}{2}(E_{\text{out}}(g_1) + E_{\text{out}}(g_2)) &= \frac{1}{2}(\mathbb{E}[(f(\mathbf{x}) - g_1(\mathbf{x}))^2] + \mathbb{E}[(f(\mathbf{x}) - g_2(\mathbf{x}))^2]) \\
&= \frac{1}{2}\mathbb{E}[(f(\mathbf{x}) - g_1(\mathbf{x}))^2 + (f(\mathbf{x}) - g_2(\mathbf{x}))^2] \\
&= \frac{1}{2}\mathbb{E}[(f(\mathbf{x})^2 - 2f(\mathbf{x})g_1(\mathbf{x}) + g_1(\mathbf{x})^2) + (f(\mathbf{x})^2 - 2f(\mathbf{x})g_2(\mathbf{x}) + g_2(\mathbf{x})^2)] \\
&= \frac{1}{2}\mathbb{E}[2f(\mathbf{x})^2 - 2f(\mathbf{x})(g_1(\mathbf{x}) + g_2(\mathbf{x})) + g_1(\mathbf{x})^2 + g_2(\mathbf{x})^2] \\
&= \mathbb{E}[f(\mathbf{x})^2 - f(\mathbf{x})(g_1(\mathbf{x}) + g_2(\mathbf{x})) + \frac{1}{2}g_1(\mathbf{x})^2 + \frac{1}{2}g_2(\mathbf{x})^2] \\
&= \mathbb{E}[f(\mathbf{x})^2 - 2f(\mathbf{x})\frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x})) + \frac{1}{2}(g_1(\mathbf{x})^2 + g_2(\mathbf{x})^2)] \\
&= \mathbb{E}[f(\mathbf{x})^2 - 2f(\mathbf{x})\frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x})) + \frac{1}{2}(g_1(\mathbf{x})^2 + 2g_1(\mathbf{x})g_2(\mathbf{x}) + g_2(\mathbf{x})^2) - g_1(\mathbf{x})g_2(\mathbf{x})] \\
&= \mathbb{E}[f(\mathbf{x})^2 - 2f(\mathbf{x})\frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x})) + \frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 + \frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 - g_1(\mathbf{x})g_2(\mathbf{x})] \\
&= \mathbb{E}[f(\mathbf{x})^2 - 2f(\mathbf{x})g(\mathbf{x}) + g(\mathbf{x})^2 + \frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 - g_1(\mathbf{x})g_2(\mathbf{x})] \\
&= \mathbb{E}[(f(\mathbf{x}) - g(\mathbf{x}))^2 + \frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 - g_1(\mathbf{x})g_2(\mathbf{x})] \\
&= \mathbb{E}[(f(\mathbf{x}) - g(\mathbf{x}))^2] + \mathbb{E}\left[\frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 - g_1(\mathbf{x})g_2(\mathbf{x})\right] \\
&= E_{\text{out}}(g) + \mathbb{E}\left[\frac{1}{4}(g_1(\mathbf{x}) + g_2(\mathbf{x}))^2 - g_1(\mathbf{x})g_2(\mathbf{x})\right] \\
&= E_{\text{out}}(g) + \mathbb{E}\left[\frac{1}{4}(g_1(\mathbf{x}) - g_2(\mathbf{x}))^2\right]
\end{aligned}$$

We use the linearity of expected values in the above calculation; $\mathbb{E}[u] + \mathbb{E}[v] = \mathbb{E}[u + v]$. Now, it suffices to show that $\mathbb{E}\left[\frac{1}{4}(g_1(\mathbf{x}) - g_2(\mathbf{x}))^2\right] \geq 0$.

$$\begin{aligned}
\mathbb{E}\left[\frac{1}{4}(g_1(\mathbf{x}) - g_2(\mathbf{x}))^2\right] &\geq 0 \\
(g_1(\mathbf{x}) - g_2(\mathbf{x}))^2 &\geq 0 \\
b = g_1(\mathbf{x}) - g_2(\mathbf{x}) \\
b^2 &\geq 0
\end{aligned}$$

Thus, we have that $E_{\text{out}}(g)$ cannot be worse than the average of $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$.

[d] *False*: $E_{\text{out}}(g)$ does not have to fall between the inclusive interval $[E_{\text{out}}(g_1), E_{\text{out}}(g_2)]$. For example if g_1 is consistently overestimating the target $y = f(\mathbf{x})$ for $\mathbf{x} \in \chi$, and g_2 is consistently underestimating the target y , then taking the average by $g(\mathbf{x}) = \frac{1}{2}(g_1(\mathbf{x}) + g_2(\mathbf{x}))$ may result in an $E_{\text{out}}(g)$ that is lower than that of both $E_{\text{out}}(g_1)$ and $E_{\text{out}}(g_2)$.