

For myself:

$$\mathbb{E}_{Y|X} (Y - c)^2 = \mathbb{E}_{Y|X} (Y^2 - 2cY + c^2) \quad (1)$$

$$= \mathbb{E}_{Y|X} Y^2 - 2c \mathbb{E}_{Y|X} Y + c^2 \quad (2)$$

$$= \left[\mathbb{E}_{Y|X} Y^2 - (\mathbb{E}_{Y|X} Y)^2 \right] + (\mathbb{E}_{Y|X} Y - c)^2, \quad (3)$$

whence it follows that

$$\arg \min \mathbb{E}_{Y|X} (Y - c)^2 = \mathbb{E}_{Y|X} Y, \quad (4)$$

$$\min \mathbb{E}_{Y|X} (Y - c)^2 = \mathbb{E}_{Y|X} Y^2 - (\mathbb{E}_{Y|X} Y)^2 \quad (5)$$

$$= \text{Var}(Y|X). \quad (6)$$

Ex. 2.1: Suppose each of K -classes has an associated target t_k , which is a vector of all zeros, except a one in the k th position. Show that classifying to the largest element of \hat{y} amounts to choosing the closest target, $\min_k \|t_k - \hat{y}\|$, if the elements of \hat{y} sum to one.

Answer:

Let $\hat{y} \in \mathbb{R}^p : \sum_i \hat{y}_i = 1$. Then, for any j ,

$$\|t_j - \hat{y}\|^2 = \sum_{i \neq j} \hat{y}_i^2 + (1 - \hat{y}_j)^2 \quad (7)$$

$$= \sum_{i \neq j} \hat{y}_i^2 + \left(\sum_{i \neq j} \hat{y}_i \right)^2 \quad (8)$$

$$= 2 \sum_{h, i \neq j; h \leq i} \hat{y}_h \hat{y}_i \quad (9)$$

$$= 2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \sum_{h=j \vee i=j, h \leq i} \hat{y}_h \hat{y}_i \quad (10)$$

$$= 2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \sum_i \hat{y}_k \hat{y}_i \quad (11)$$

$$= 2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \hat{y}_k \sum_i \hat{y}_i \quad (12)$$

$$= 2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \hat{y}_k \quad (13)$$

$$(14)$$

So, if the closest thingy is the k th,

$$\|t_k - \hat{y}\| \leq \|t_j - \hat{y}\| \forall j, \quad (15)$$

$$\|t_k - \hat{y}\|^2 \leq \|t_j - \hat{y}\|^2 \forall j, \quad (16)$$

$$2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \hat{y}_k \leq 2 \sum_{h \leq i} \hat{y}_h \hat{y}_i - 2 \hat{y}_j \forall j, \quad (17)$$

$$\hat{y}_k \geq \hat{y}_j \forall j. \quad (18)$$

Ex. 2.2: Show how to compute the Bayes decision boundary for the simulation example in Figure 2.5.

Answer:

Maybe the previous exercise was meant to help with this one, since we have to pick k such that $p_k \geq p_j \forall j$ with p in the probability simplex, which is equivalent to picking the axis closest to p . But that doesn't look less

computationally expensive. In the general case, I guess you could construct a grid and simply check for those conditions.

In the case of a binary classification problem, you can take advantage in the fact that the decision boundary takes the form $p_1 = 0.5$, and just move along the unique choice perpendicular to the gradient of p_1 .

Ex. 2.3: Derive equation (2.24).

Answer:

We have a uniform probability distribution along $S_p = \{x \in \mathbb{R}^p : \|x\| \leq 1\}$, and we take N samples. The median minimum distance is such that the probability of a smaller minimum distance be $1/2$. The probability that the minimum distance is above r is the probability that all N samples are above r , which should be the N th power of the probability that any individual sample is above r , which is $1 - r^p$. That is,

$$P(\|x\| \geq r) = (1 - r^p)^N \quad (19)$$

and, setting this to equal $1/2$, we get

$$d(p, N) = \left(1 - \frac{1}{2^{1/N}}\right)^{1/p}. \quad (20)$$

Thank God for commulative probabilities.

Ex. 2.4: The edge effect problem discussed on page 15 is not peculiar to uniform sampling from bounded domains. Consider inputs drawn from a spherical multinormal distribution $X \sim N(0, \mathbf{I}_p)$. The squared distance from any sample point to the origin has a χ_p^2 distribution with mean p . Consider a prediction point x_0 drawn from this distribution, and let $a = x_0/\|x_0\|$ be an associated unit vector. Let $z_i = a^T x_i$ be the projection of each of the training points in this direction.

- a Show that the z_i are distributed $N(0, 1)$ with expected squared distance from the origin 1, while the target point has expected squared distance p from the origin.
- b For $p = 10$ show that the expected distance of a test point from the center of the training data is 3.1 standard deviations, while all the training points have expected distance 1 along direction a . So most prediction points see themselves as lying on the edge of the training set.

Answer:

Part a:

$N(0, \mathbf{I}_p)$ is precisely obtained as the probability distribution of a point each of whose components are taken from $N(0, 1)$, which shows $a^T x_i$ is taken from a distribution $N(0, 1)$.

Long road: Let $P_X = N(\mu, \Sigma)$, which is to say

$$P_X(x) = \frac{1}{\sqrt{2\pi \det \Sigma}^p} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right). \quad (21)$$

Then,

$$P_{a^T X}(y) = \int_{\mathbb{R}^p} dx P_X(x) \delta(a^T x - y) \quad (22)$$

$$= \int_{\mathbb{R}^p: a^T x = y} dx P_X(x) \frac{1}{\|a\|} \cdot 1 \quad (23)$$

We now take into account that $a^T x = y$ means $x = \sum_i x_i u_i + ya = x_\perp + ya$ for $\{a, u_i\}$ an orthonormal basis. This way, and decomposing $\mu = \mu_\perp + \mu_a a$,

$$(x - \mu)^T \Sigma^{-1}(x - \mu) = (x_\perp - \mu_\perp + (y - \mu_a)a)^T \Sigma^{-1}(x_\perp - \mu_\perp + (y - \mu_a)a) \quad (24)$$

$$= (x_\perp - \mu_\perp)^T \Sigma^{-1}(x_\perp - \mu_\perp) \quad (25)$$

$$+ 2(x_\perp - \mu_\perp)^T \Sigma^{-1}(y - \mu_a)a \quad (26)$$

$$+ (y - \mu_a)^2 a^T \Sigma^{-1}a. \quad (27)$$

And so we plug all that into the big integral. The (exponential of the) last term can be taken outside the integral, whereas the other two represent stuff to integrate. The contributions from $x_\perp - \mu_\perp$ and $\mu_\perp - x_\perp$ cancel out the second factor, and so it ends up being a constant factor. We can then conclude the probability distribution for y is $N(\mu_a, a^T \Sigma^{-1} a)$.

On the other hand, the expected square distance $E\|x\|^2$ is given by

$$E\|x\|^2 = \frac{\int dx \|x\|^2 \exp\left(-\frac{1}{2}\|x\|^2\right)}{\int dx \exp\left(-\frac{1}{2}\|x\|^2\right)} \quad (28)$$

$$= \frac{\int_0^\infty dr r^{p-1} \int d\Theta r^2 \exp\left(-\frac{1}{2}r^2\right)}{\int_0^\infty dr r^{p-1} \int d\Theta \exp\left(-\frac{1}{2}r^2\right)} \quad (29)$$

$$= \frac{\int_0^\infty dr r^{p+1} \exp\left(-\frac{1}{2}r^2\right)}{\int_0^\infty dr r^{p-1} \exp\left(-\frac{1}{2}r^2\right)} \quad (30)$$

$$= \frac{\int_0^\infty d\left(\frac{1}{2}r^2\right) 2^{\frac{1}{2}p} \left(\frac{1}{2}r^2\right)^{\frac{1}{2}p} \exp\left(-\frac{1}{2}r^2\right)}{\int_0^\infty d\left(\frac{1}{2}r^2\right) 2^{\frac{1}{2}p-1} \left(\frac{1}{2}r^2\right)^{\frac{1}{2}p-1} \exp\left(-\frac{1}{2}r^2\right)} \quad (31)$$

$$= 2 \frac{\int_0^\infty dt t^{\frac{1}{2}p} \exp(-t)}{\int_0^\infty dt t^{\frac{1}{2}p-1} \exp(-t)} \quad (32)$$

$$= 2 \frac{\Gamma\left(\frac{1}{2}p + 1\right)}{\Gamma\left(\frac{1}{2}p\right)} \quad (33)$$

$$= p. \quad (34)$$

Note to self: Gamma thingy can be done through parts:

$$\Gamma(x+1) = \int_0^\infty dt t^x \exp(-t) \quad (35)$$

$$= \left[-xt^{x-1} \exp(-t) \right]_0^\infty + \int_0^\infty dt xt^{x-1} \exp(-t) \quad (36)$$

$$= x\Gamma(x) \quad (37)$$

Part b:

If we compute the expected distance as $\sqrt{E\|x\|^2}$ we're left with about 3.16, which I guess is what the author did. And here I was, about to compute $E\|x\|$. In fact, let's do it for the lulz.

$$E\|x\| = \frac{\int dr \|x\| \exp\left(-\frac{1}{2}\|x\|^2\right)}{\int dx \exp\left(-\frac{1}{2}\|x\|^2\right)} \quad (38)$$

$$= \dots \quad (39)$$

$$= \frac{\int dr r^p \exp\left(-\frac{1}{2}r^2\right)}{\int dx r^{p-1} \exp\left(-\frac{1}{2}r^2\right)} \quad (40)$$

$$= \sqrt{2} \frac{\Gamma\left(\frac{1}{2}p + \frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}p\right)} \quad (41)$$

$$\approx 3.08. \quad (42)$$

For the expected distance along direction a , we get 1 as the square root of the expected square distance, and

$$\mathbb{E}\|x\| = \frac{\int_0^\infty dr r \exp(-\frac{1}{2}r^2)}{\int_0^\infty dr \exp(-\frac{1}{2}r^2)} \quad (43)$$

$$= \frac{\int_0^\infty dt \exp(-t)}{\sqrt{\frac{\pi}{2}}} \quad (44)$$

$$= \frac{2}{\pi} \quad (45)$$

$$\approx 0.78. \quad (46)$$

Ex. 2.5:

- a Derive equation (2.27). The last line makes use of (3.8) through a conditioning argument.
- b Derive equation (2.28), making use of the *cyclic* property of the trace operator $[\text{Tr}(AB) = \text{Tr}(BA)]$, and its linearity (which allows to interchange the order of trace and expectation).

Answer:

The relationship between the random variables X and Y is

$$Y = X^T \beta + \varepsilon, \quad (47)$$

with $\varepsilon \sim N(0, \sigma)$ and X taking values in \mathbb{R}^p . As a training set, we have N inputs arranged in row form into matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$, and N outputs arranged in $\mathbf{y} \in \mathbb{R}^{N \times 1}$ of the form $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon}$ is a vector of samples from ε 's probability distribution function.

So, the expected prediction error at input value x_0 is given by

$$\text{EPE}(x_0) = \mathbb{E}_{y_0|x_0} \mathbb{E}_{\mathcal{T}} (y_0 - \hat{y}_0)^2, \quad (48)$$

$$\mathbb{E}_{y_0|x_0} = \int dy_0 P(y_0|x_0), \quad (49)$$

$$\mathbb{E}_{\mathcal{T}} = \mathbb{E}_{\mathbf{X}} \mathbb{E}_{\varepsilon} \quad (50)$$

$$= \int d\mathbf{X} P(\mathbf{X}) \int d\boldsymbol{\varepsilon} P(\boldsymbol{\varepsilon}). \quad (51)$$

We observe that all y_0 (ε), \mathbf{X} and $\boldsymbol{\varepsilon}$ are independent variables, whence the expectation values commute.

The estimated value is defined by

$$\hat{y}_0 = x_0^T \hat{\beta}, \quad (52)$$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (53)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \boldsymbol{\varepsilon}) \quad (54)$$

$$= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}. \quad (55)$$

We thus have:

$$\mathbb{E}_{y_0|x_0} y_0 = \mathbb{E}_{\varepsilon|x_0} (x_0^T \beta + \varepsilon) \quad (56)$$

$$= x_0^T \beta + \cancel{\mathbb{E}_{\varepsilon|x_0} \varepsilon}, \quad 0 \quad (57)$$

$$\text{Var}(y_0|x_0) = \mathbb{E}_{\varepsilon|x_0} (y_0 - \mathbb{E}_{\varepsilon|x_0} y_0)^2 \quad (58)$$

$$= \mathbb{E}_{\varepsilon|x_0} (x_0^T \beta + \varepsilon - x_0^T \beta)^2 \quad (59)$$

$$= \mathbb{E}_{\varepsilon|x_0} \varepsilon^2 = 1, \quad (60)$$

$$\mathbb{E}_{\mathcal{T}} \hat{y}_0 = \mathbb{E}_{\mathcal{T}} (x_0^T \hat{\beta}) \quad (61)$$

$$= \mathbb{E}_{\mathcal{T}} (x_0^T \beta + \varepsilon) \quad (62)$$

$$= x_0^T \beta + \cancel{\mathbb{E}_{\mathcal{T}} \varepsilon}, \quad 0 \quad (63)$$

$$\text{Var} \hat{y}_0 = \mathbb{E}_{\mathcal{T}} (\hat{y}_0 - \mathbb{E}_{\mathcal{T}} \hat{y}_0)^2 \quad (64)$$

$$= \mathbb{E}_{\mathcal{T}} \left(x_0^T \beta + x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} - x_0^T \beta \right)^2 \quad (65)$$

$$= \mathbb{E}_{\mathcal{T}} \left(x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \right)^2 \quad (66)$$

$$= \mathbb{E}_{\mathcal{T}} x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x_0 \quad (67)$$

$$= \mathbb{E}_{\mathbf{X}} x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}_{\boldsymbol{\varepsilon}} (\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x_0 \quad (68)$$

$$= \mathbb{E}_{\mathbf{X}} x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_N \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x_0 \quad (69)$$

$$= \sigma^2 \mathbb{E}_{\mathbf{X}} x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0 \quad (70)$$

$$(71)$$

Now,

$$\text{EPE}(x_0) = \mathbb{E}_{y_0|x_0} \mathbb{E}_{\mathcal{T}} (y_0 - \hat{y}_0)^2 \quad (72)$$

$$= \mathbb{E}_{y_0|x_0} y_0^2 - 2\mathbb{E}_{y_0|x_0} y_0 \mathbb{E}_{\mathcal{T}} \hat{y}_0 + \mathbb{E}_{\mathcal{T}} \hat{y}_0^2 \quad (73)$$

$$= \text{Var}(y_0|x_0) + (\mathbb{E}_{y_0|x_0} y_0)^2 - 2\mathbb{E}_{y_0|x_0} y_0 \mathbb{E}_{\mathcal{T}} \hat{y}_0 + \text{Var}_{\mathcal{T}}(\hat{y}_0) + (\mathbb{E}_{\mathcal{T}} \hat{y}_0)^2 \quad (74)$$

$$= \text{Var}(y_0|x_0) + \text{Var}_{\mathcal{T}}(\hat{y}_0) + (\mathbb{E}_{y_0|x_0} y_0 - \mathbb{E}_{\mathcal{T}} \hat{y}_0)^2 \quad (75)$$

$$= \text{Var}(y_0|x_0) + \text{Var}_{\mathcal{T}}(\hat{y}_0) + \text{Bias}^2(\hat{y}_0) \quad (76)$$

$$= \sigma^2 + \sigma^2 x_0^T \mathbb{E}_{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} x_0 + 0. \quad (77)$$

Since

$$\mathbf{X}^T \mathbf{X} = \left(\sum_{\mu} x_i^{\mu} x_j^{\mu} \right), \quad (78)$$

and assuming $EX = 0$ for each element, this is approximately $NCov(X)$. So,

$$E_{x_0} EPE(x_0) = \sigma^2 + \frac{\sigma^2}{N} E_{x_0} (x_0^T Cov^{-1}(X) x_0) \quad (79)$$

$$= \sigma^2 + \frac{\sigma^2}{N} E_{x_0} Tr (x_0^T Cov^{-1}(X) x_0) \quad (80)$$

$$= \sigma^2 + \frac{\sigma^2}{N} Tr (Cov^{-1}(X) E_{x_0} x_0 x_0^T) \quad (81)$$

$$= \sigma^2 + \frac{\sigma^2}{N} Tr (Cov^{-1}(X) Cov(x_0)) \quad (82)$$

$$= \sigma^2 + \frac{\sigma^2}{N} Tr (I_p) \quad (83)$$

$$= \sigma^2 + \sigma^2 \frac{p}{N} \quad (84)$$

$$(85)$$

Ex. 2.6: Consider a regression problem with inputs x_i and outputs y_i , and a parametrized model $f_\theta(x)$ to be fit by least squares. Show that if there are observations with *tied* or *identical* values of x , then the fit can be obtained from a reduced weighted least squares problem.

Answer:

Let $I_{i,j}$ be the index for the i th unique value X_i (from 1 to U), with j ranging from 1 to the number of repetitions R_i . Then,

$$RSS(\theta) = \sum_{i=1}^N (y_i - f_\theta(x_i))^2 \quad (86)$$

$$= \sum_{i=1}^U \sum_{j=1}^{R_i} (y_{I_{i,j}} - f_\theta(x_{I_{i,j}}))^2 \quad (87)$$

$$= \sum_{i=1}^U \sum_{j=1}^{R_i} (y_{I_{i,j}} - f_\theta(X_i))^2 \quad (88)$$

$$= \sum_{i=1}^U \sum_{j=1}^{R_i} (y_{I_{i,j}}^2 - 2y_{I_{i,j}} f_\theta(X_i) + f_\theta^2(X_i)) \quad (89)$$

$$= \sum_{i=1}^U \sum_{j=1}^{R_i} (y_{I_{i,j}}^2 - 2y_{I_{i,j}} f_\theta(X_i) + f_\theta^2(X_i)) \quad (90)$$

$$= \sum_{i=1}^U R_i (-2\bar{y}_i f_\theta(X_i) + f_\theta^2(X_i)) + \sum_{i=1}^U \sum_{j=1}^{R_i} y_{I_{i,j}}^2 \quad (91)$$

$$= \sum_{i=1}^U R_i (\bar{y}_i - f_\theta(X_i))^2 - \sum_{i=1}^U R_i \bar{y}_i^2 + \sum_{i=1}^U \sum_{j=1}^{R_i} y_{I_{i,j}}^2, \quad (92)$$

$$\bar{y}_i = \frac{1}{R_i} \sum_{j=1}^{R_i} y_{ij}. \quad (93)$$

Since the later two sumands are constant with respect to θ , the minimum of this function equals the minimum of

$$RSS(\theta) = \sum_{i=1}^U R_i (\bar{y}_i - f_\theta(X_i))^2. \quad (94)$$

Ex. 2.7: Suppose we have a sample of N pairs x_i, y_i drawn i.i.d. from the distribution characterized as follows:

- $x_i \sim h(x)$, the design density
- $y_i = f(x_i) + \varepsilon_i$, f is the regression function
- $\varepsilon_i \sim (0, \sigma^2)$ (mean zero, variance σ^2)

We construct an estimator for f linear in the y_i ,

$$\hat{f}(x_0) = \sum_{i=1}^N l_i(x_0; \mathcal{X}) y_i \quad (95)$$

where the weights $l_i(x_0; \mathcal{X})$ do not depend on the y_i , but do depend on the entire training sequence of x_i , denoted here by \mathcal{X} .

- Show that linear regression and k -nearest-neighbor regression are members of this class of estimators. Describe explicitly the weights $l_i(x_0; \mathcal{X})$ for both these cases.
- Decompose the conditional mean-squared error

$$\mathbb{E}_{\mathcal{Y}|\mathcal{X}} \left(f(x_0) - \hat{f}(x_0) \right)^2 \quad (96)$$

into a conditional squared bias and a conditional variance component. Like \mathcal{X} , \mathcal{Y} represents the entire training sequence of y_i .

- Decompose the (unconditional) mean-squared error

$$\mathbb{E}_{\mathcal{X}, \mathcal{Y}} \left(f(x_0) - \hat{f}(x_0) \right)^2 \quad (97)$$

into a squared bias and a variance component.

- Establish a relationship between the squared biases and variances in the above two cases.

Answer:

Part a:

Linear regression prescribes

$$\hat{f}(x_0) = \mathbf{l}^T(x_0; \mathcal{X}) \mathbf{y} \quad (98)$$

$$= \sum_i^N l_i(x_0; \mathcal{X}) y_i, \quad (99)$$

$$\mathbf{l}(x_0; \mathcal{X}) = x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T. \quad (100)$$

KNN prescribes

$$\hat{f}(x_0) = \sum_{i=1}^N l_i(x_0; \mathcal{X}) y_i, \quad (101)$$

$$l_i(x_0; \mathcal{X}) = \frac{1}{K} \chi_{\text{KNN}(x_0)}(x_i), \quad (102)$$

$$\text{KNN}(x_0) : \text{ set of } K \text{ nearest neighbors to } x_0. \quad (103)$$

Part b and c:

$$\mathbb{E} \left(f(x_0) - \hat{f}(x_0) \right)^2 = f^2(x_0) - 2f(x_0)\mathbb{E}\hat{f}(x_0) + \mathbb{E}\hat{f}^2(x_0) \quad (104)$$

$$= \left(f(x_0) - \mathbb{E}\hat{f}(x_0) \right)^2 + \mathbb{E}\hat{f}^2(x_0) - (\mathbb{E}f(x_0))^2 \quad (105)$$

$$= \text{Bias}^2 \left(\hat{f}(x_0) \right) + \text{Var} \left(\hat{f}(x_0) \right), \quad (106)$$

since in both cases x_0 and $f(x_0)$ are constant.

Part d:

In part b, the sample set is constant, whereas in part c it is averaged. But we can decompose $E_{\mathcal{X}, \mathcal{Y}}$ as $E_{\mathcal{X}} E_{\mathcal{Y}|\mathcal{X}}$. So, for the bias, we have

$$\text{Bias}_{\mathcal{X}, \mathcal{Y}}(\hat{f}(x_0)) = E_{\mathcal{X}, \mathcal{Y}} \hat{f}(x_0) - f(x_0) \quad (107)$$

$$= E_{\mathcal{X}} E_{\mathcal{Y}|\mathcal{X}}(\hat{f}(x_0) - f(x_0)) \quad (108)$$

$$= E_{\mathcal{X}} \text{Bias}_{\mathcal{Y}|\mathcal{X}}(\hat{f}(x_0)). \quad (109)$$

Now, in this particular case,

$$\text{Bias}_{\mathcal{Y}|\mathcal{X}}(\hat{f}(x_0)) = E_{\mathcal{Y}|\mathcal{X}}(\hat{f}(x_0) - f(x_0)) \quad (110)$$

$$= E_{\mathcal{Y}|\mathcal{X}}\left(\sum_i l_i(x_0|\mathcal{X}) y_i - f(x_0)\right) \quad (111)$$

$$= E_{\mathcal{Y}|\mathcal{X}}\left(\sum_i l_i(x_0|\mathcal{X})(f(x_i) + \varepsilon_i) - f(x_0)\right) \quad (112)$$

$$= \sum_i l_i(x_0|\mathcal{X}) f(x_i) + E_{\varepsilon} \sum_i l_i(x_0|\mathcal{X}) \varepsilon_i - f(x_0) \quad (113)$$

$$= \sum_i l_i(x_0|\mathcal{X}) f(x_i) - f(x_0). \quad (114)$$

So, the conditional bias cannot be made null without knowledge of $f(x_0)$, since it depends on the particular samples drawn. However, KNN is unbiased when averaged over training samples. To see this, we note the leftmost term is simply the sum over the K nearest neighbors. For the region where the K nearest neighbors are the ones corresponding to the samples with index I , the indexes may be relabeled so that $I = [1, K]_{\mathbb{N}}$. There are K taken from N such regions. Of these, for the region where the i th index corresponds to the K th closest neighbor, again the indexes may be relabeled so that the first one is the closest one. There are N such regions. So we have

$$E_{\mathcal{X}} \sum_{i=1}^N l_i(x_0|\mathcal{X}) f(x_i) = \prod_{i=1}^N \int_{-\infty}^{\infty} dx_i h(x_i - x_0) \sum_{j=1}^N l_j(x_0|\mathcal{X}) f(x_j - x_0) \quad (115)$$

$$= N \binom{N}{K} \int_0^{\infty} dx_1 (h(x_1) + h(-x_1)) \prod_{i=2}^K \int_{-x_1}^{x_1} dx_i h(x_i) \frac{1}{K} \sum_j f(x_j + x_0) \left(\int_{-\infty}^{-x_1} dx h(x) + \int_{x_1}^{\infty} dx h(x) \right) \quad (116)$$

Aaand I'm at a loss here.

Ex. 2.8: Compare the classification performance of linear regression and k -nearest neighbor classification on the zipcode data. In particular, consider only the 2's and 3's, and $k = 1, 3, 5, 7$ and 15. Show both the training and test error for each choice. The zipcode data are available from the book website www-stat.stanford.edu/ElemStatLearn.

Answer:

The code “2.8.py” produces figure 1. It's noticeable that linear regression performs worse with the test data, but generalizes better, in this case, than all other instances of KNN.

Ex. 2.9: Consider a linear regression model with p parameters, fit by least squares to a set of training data $(x_1, y_1), \dots, (x_N, y_N)$ drawn at random from a population. Let $\hat{\beta}$ be the least squares estimate. Suppose we have some test data $(\tilde{x}_1, \tilde{y}_1), \dots, (\tilde{x}_M, \tilde{y}_M)$ drawn at random from the same population as the training data. If $R_{tr}(\beta) = \frac{1}{N} \sum_i^N (y_i - \beta^T x_i)^2$ and $R_{te}(\beta) = \frac{1}{M} \sum_i^M (\tilde{y}_i - \beta^T \tilde{x}_i)^2$, prove that

$$E[R_{tr}(\hat{\beta})] \leq E[R_{te}(\hat{\beta})], \quad (117)$$

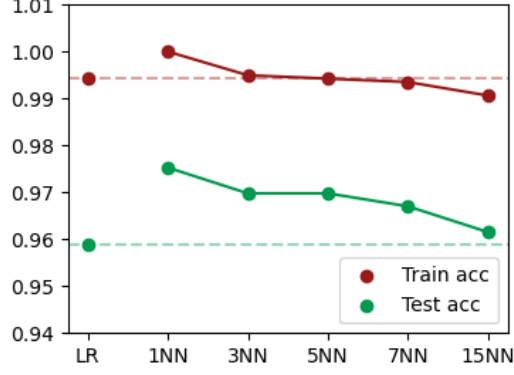


Figure 1: Accuracy as a function of the regressor. “LR” corresponds to Linear Regressor, and “ i NN” corresponds to an i -nearest neighbors regressor.

where the expectations are over all that is random in each expression.

Answer:

Bear in mind that $Ef(x_i, y_i, \hat{\beta}) = Ef(x_j, y_j, \hat{\beta})$. To see this, first note $\hat{\beta}$ is invariant with respect to changes of labels:

$$\hat{\beta}(\mathbf{X}, \mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (118)$$

$$= (\mathbf{X}^T S_{ij}^T S_{ij} \mathbf{X})^{-1} \mathbf{X}^T S_{ij}^T S_{ij} \mathbf{y} \quad (119)$$

$$= \left((S_{ij} \mathbf{X})^T (S_{ij} \mathbf{X}) \right)^{-1} (S_{ij} \mathbf{X})^T (S_{ij} \mathbf{y}) \quad (120)$$

$$= \hat{\beta}(S_{ij} \mathbf{X}, S_{ij} \mathbf{y}), \quad (121)$$

where $S_{ij} = S_{ij}^T = S_{ij}^{-1}$ is the i to j row operator. It then follows that

$$Ef(x_i, y_i, \hat{\beta}) = \int P(dx_1, dy_1) \dots P(dx_N, dy_N) f(x_i, y_i, \hat{\beta}) \quad (122)$$

$$= \int P(dx_1, dy_1) \dots P(dx_N, dy_N) f(x_j, y_j, \hat{\beta}) \quad (\text{Change of labels}) \quad (123)$$

$$= Ef(x_j, y_j, \hat{\beta}). \quad (124)$$

Thus,

$$E[R_{tr}(\hat{\beta}(\mathbf{X}, \mathbf{y}))] = E \left[\frac{1}{N} \sum_i^N \left(y_i - \hat{\beta}^T(\mathbf{X}, \mathbf{y}) x_i \right)^2 \right] \quad (125)$$

$$= \frac{1}{N} \sum_i^N E \left(y_i - \hat{\beta}^T(\mathbf{X}, \mathbf{y}) x_i \right)^2 \quad (126)$$

$$= E \left(y_1 - \hat{\beta}^T(\mathbf{X}, \mathbf{y}) x_1 \right)^2, \quad (127)$$

$$E[R_{ts}(\hat{\beta}(\mathbf{X}, \mathbf{y}))] = E \left(\tilde{y}_1 - \hat{\beta}^T(\mathbf{X}, \mathbf{y}) \tilde{x}_1 \right)^2. \quad (128)$$

Now, for the comparison,

$$E[R_{tr}(\hat{\beta}(\mathbf{X}, \mathbf{y}))] = E \left(y_1 - \hat{\beta}^T(\mathbf{X}, \mathbf{y}) x_1 \right)^2 \quad (129)$$

$$= Ey_1^2 - 2Ey_1 \hat{\beta}^T(\mathbf{X}, \mathbf{y}) x_1 + E \left[\hat{\beta}(\mathbf{X}, \mathbf{y})^T x_1 \right]^2 \quad (130)$$

Again I'm at a loss. It's time to move on.