

1.

$$\|t_k - \hat{y}\| = \sqrt{\sum_{i \neq k} \hat{y}_i^2 + (1 - 2\hat{y}_k + \hat{y}_k^2)} = \sqrt{\sum_i \hat{y}_i^2 + (1 - 2\hat{y}_k)}.$$

In the right hand side, the first term  $\sum_i \hat{y}_i^2$  is independent of  $k$ , so we need only consider the second term, which is minimized when  $\hat{y}_k$  is maximized. Thus,  $\operatorname{argmin} \|t_k - \hat{y}\|$  is equal to the  $k$  such that  $\hat{y}_k$  is maximized.

2. Let the orange means have corresponding pdfs  $o_i$  and the blue means have corresponding pdfs  $b_i$ . Then for any individual mean  $m_*$ , the probability that a particular point  $p = (x_0, y_0)$  was sampled from the corresponding distribution  $D_*$  is  $P(D_*|p) = P(p|D_*)P(D_*)/(\sum_i P(p|D_i)P(D_i))$ . Since the means are sampled with equal probabilities and there are an equal number of blue and orange points, all the prior probabilities  $P(D_i)$  are equal so the above reduces to  $P(D_*|p) = P(p|D_*)/(\sum_i P(p|D_i)) = f_*(p)/\sum_i f_i(p)$ . Summing over all the orange means, the probability that a point  $p$  is orange is  $\sum_i o_i(p)/\sum_i o_i(p) + b_i(p)$  and we classify  $p$  as orange if this probability is larger than  $1/2$ . This yields the Bayes decision boundary where this probability is equal to  $1/2$  (or equivalently  $\sum_i o_i(p) = \sum_i b_i(p)$ ).
3. For a uniformly distributed point in the  $p$ -dimensional unit sphere, the cdf of the distance to the origin is  $P(R < r) = cr^p/c1^p = r^p$ . For  $N$  uniformly (independently) distributed points in the  $p$ -dimensional unit sphere, the minimum distance thus has cdf  $P(\min(R_i) < r) = 1 - P(\min(R_i) \geq r) = 1 - P(R_1 \geq r)^N = 1 - (1 - r^p)^N$ . The median of this distribution occurs when  $P(\min(R_i) < r) = 1/2$ , so solving  $1/2 = 1 - (1 - r^p)^N$  for  $r$  yields  $r = (1 - (1/2)^{1/N})^{1/p}$ .
4. Since  $a$  is a unit vector,  $\|a\|^2 = 1$ . Since  $x_i \sim N(0, I_p)$ , each component  $(x_i)_j$  is an independent  $N(0, 1)$  random variable and thus by properties of the normal distribution,  $a^T x_i \sim N(\sum_i a_i 0 = 0, \sum_i a_i^2 1) = N(0, \|a\|^2) = N(0, 1)$ . Thus  $E(z_i^2) = \operatorname{Var}(z_i) = 1$ .
5. (a) In the prediction situation where the true relationship  $Y = X^T \beta + \varepsilon$  is linear (up to noise  $\varepsilon \sim N(0, \sigma^2)$ ), there are two sources of error: first,  $y_0$  is not deterministically determined by  $x_0$  so there is error coming from the model itself and second, the estimation of  $\beta$  depends on the training data we see.

Thus,  $EPE(x_0) = E((y_0 - \hat{y}_0)^2)$  where the expectation is over both  $y_0$  (given  $x_0$ ) and the training data  $T$ . (Note that although  $X$  is typically considered fixed, we can also view the training data  $T$  as randomly sampled pairs  $(X, Y)$  from some underlying joint distribution. Furthermore,  $y_0|x_0$  and  $T$  are independent since the training and testing data are independent, so the joint expectation can be written successively as  $E_{y_0|x_0} E_T$ .) Adding and subtracting  $x_0^T \beta$  yields  $E((y_0 - x_0^T \beta)^2 + 2(y_0 - x_0^T \beta)(x_0^T \beta - \hat{y}_0) + (x_0^T \beta - \hat{y}_0)^2)$ . This further simplifies to  $E(\varepsilon_0) + 2E(\varepsilon_0(x_0^T \beta - \hat{y}_0)) + E((x_0^T \beta - \hat{y}_0)^2)$ . The first term does not depend on  $T$ , so the expectation over  $T$  is trivial and  $E(\varepsilon_0^2) = \sigma^2$  by definition of the error term. (More precisely,  $E((y_0 - x_0^T \beta)^2) = E_{y_0|x_0}((y_0 - x_0^T \beta)^2) = \operatorname{Var}_{y_0|x_0}(y_0 - x_0^T \beta) + E_{y_0|x_0}(y_0 - x_0^T \beta)^2 = \operatorname{Var}_{y_0|x_0}(y_0) = \sigma^2$ .) Next, since the error term  $\varepsilon_0 = y_0 - x_0^T \beta$  does not depend on  $T$  while the expectation  $x_0^T \beta - \hat{y}_0$  does not depend on  $y_0$ ,

the second term factors as  $E_{y_0|x_0}(\varepsilon_0)E_T(x_0^T\beta - \hat{y}_0)$ . By the model definition, both terms in the product above are zero. Thus we are left with  $\sigma^2 + E_T((x_0^T\beta - x_0^T\hat{\beta})^2)$ . Focusing on the remaining expectation, it is equal to  $Var_T(x_0^T\beta - x_0^T\hat{\beta}) + E_T(x_0^T\beta - x_0^T\hat{\beta})^2 = Var_T(x_0^T\hat{\beta}) + (x_0^T\beta - E_T(x_0^T\hat{\beta}))^2$ . Since  $\hat{\beta}$  is unbiased when the model is linear, the latter term is zero (since it is the squared bias of the estimator  $x_0^T\hat{\beta}$  for the estimand  $x_0^T\beta$ ). Simplifying the variance further:

$$\begin{aligned} Var_T(x_0^T\hat{\beta}) &= x_0^T Var_T((X^T X)^{-1} X^T Y) x_0 = x_0^T (X^T X)^{-1} X^T Var_T(Y) X (X^T X)^{-1} x_0 \\ &= x_0^T (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} x_0 = x_0^T (X^T X)^{-1} x_0 \sigma^2 \end{aligned}$$

under the assumption that the errors are uncorrelated and viewing  $X$  as fixed. If we do not view  $X$  as fixed, then

$$\begin{aligned} &Var_T((X^T X)^{-1} X^T Y) \\ &= E_T((X^T X)^{-1} X^T Y Y^T X (X^T X)^{-1}) - E_T((X^T X)^{-1} X^T Y) E_T(Y^T X (X^T X)^{-1}) \\ &= E_T((\beta + (X^T X)^{-1} X^T \varepsilon)(\beta^T + \varepsilon^T X (X^T X)^{-1})) \\ &\quad - E_T(\beta + (X^T X)^{-1} X^T \varepsilon) E_T(\beta^T + \varepsilon^T X (X^T X)^{-1}) \\ &= E_T(\beta\beta^T) + E_T((X^T X)^{-1} X^T \varepsilon\beta^T) + E_T(\beta\varepsilon^T X (X^T X)^{-1}) + E_T((X^T X)^{-1} X^T \varepsilon\varepsilon^T X (X^T X)^{-1}) \\ &\quad - E_T(\beta) E_T(\beta^T) - E_T((X^T X)^{-1} X^T \varepsilon) E_T(\beta^T) - E_T(\beta) E_T(\varepsilon^T X (X^T X)^{-1}) \\ &\quad - E_T((X^T X)^{-1} X^T \varepsilon) E_T(\varepsilon^T X (X^T X)^{-1}). \end{aligned}$$

$\beta$  is fixed (but unknown) and in particular is constant with respect to  $T$  so can be treated as a constant in all the expectations above. By Adam's law  $E_T((X^T X)^{-1} X^T \varepsilon) = E_T((X^T X)^{-1} X^T E_T(\varepsilon|X)) = 0$  since  $E(\varepsilon|X) = 0$ . After cancellation, the above becomes

$$\begin{aligned} &= E_T((X^T X)^{-1} X^T \varepsilon\varepsilon^T X (X^T X)^{-1}) = E_T((X^T X)^{-1} X^T E_T(\varepsilon\varepsilon^T|X) X (X^T X)^{-1}) \\ &= E_T((X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1}) = E_T(X^T X^{-1}) \sigma^2. \end{aligned}$$

(This result seems to be inconsistent with respect to assuming that  $X$  is fixed.)

- (b) Continuing not to assume that  $X$  is fixed, we now further assume that  $E(X) = 0$  by centering (the columns of)  $X$  if necessary. Under this assumption,  $X^T X$  converges to  $NCov(X)$  by the law of large numbers. (Otherwise, it would converge to  $NCov(X) - E(X^T)E(X)$ .) Thus,  $(X^T X)^{-1}$  converges to  $Cov(X)^{-1}/N$  (note that we do need the law of large numbers here, since even though  $E_T(X^T X) = NCov(X)$ , expectation is only linear so in particular does not necessarily interact well with inverses: it will not be true in general that as a result,  $E_T((X^T X)^{-1}) = Cov(X)^{-1}/N$ ). Applying this to the result of the previous part,

$$\begin{aligned} &E_{x_0}(\sigma^2 + x_0^T E_T((X^T X)^{-1}) x_0 \sigma^2) = \sigma^2 + \sigma^2/N E_{x_0}(x_0^T Cov(X)^{-1} x_0) \\ &= \sigma^2 + \sigma^2/N Tr(E_{x_0}(x_0 x_0^T Cov(X)^{-1})) = \sigma^2 + \sigma^2/N Tr(Cov(X)^{-1} E_{x_0}(x_0 x_0^T)) \\ &= \sigma^2 + \sigma^2/N Tr(Cov(X)^{-1} Cov(X)) = \sigma^2 + \sigma^2 p/N. \end{aligned}$$

6. Let  $x_i$  be only the distinct values of the inputs  $x_i$  and  $y_{ij}$  be all the output values corresponding to the input  $x_i$ . Then the (unweighted) RSS is  $\sum_i \sum_j (y_{ij} - f_\theta(x_i))^2$ . Let  $m_i$  be the number of times that input  $x_i$  is repeated. Then expanding the inner sum yields

$$\begin{aligned} \sum_j (y_{ij} - f_\theta(x_i))^2 &= \sum_j y_{ij}^2 - 2 \sum_j y_{ij} f_\theta(x_i) + m_i f_\theta(x_i)^2 = \sum_j y_{ij}^2 - 2 m_i \bar{y}_i f_\theta(x_i) + m_i f_\theta(x_i)^2 \\ &= \left( \sum_j y_{ij}^2 - m_i \bar{y}_i^2 \right) + m_i (\bar{y}_i^2 - 2 \bar{y}_i f_\theta(x_i) + f_\theta(x_i)^2) = \left( \sum_j y_{ij}^2 - m_i \bar{y}_i^2 \right) + m_i (\bar{y}_i - f_\theta(x_i))^2 \end{aligned}$$

and returning to the overall sum

$$\begin{aligned} \sum_i \sum_j (y_{ij} - f_\theta(x_i))^2 &= \sum_i \left( \sum_j y_{ij}^2 - m_i \bar{y}_i^2 \right) + m_i (\bar{y}_i - f_\theta(x_i))^2 \\ &= \sum_i \left( \sum_j y_{ij}^2 - m_i \bar{y}_i^2 \right) + \sum_i m_i (\bar{y}_i - f_\theta(x_i))^2. \end{aligned}$$

The second term on the RHS above is a weighted least squares problem and the first term does not depend on  $\theta$ . In particular, minimizing the RSS  $\sum_i \sum_j (y_{ij} - f_\theta(x_i))^2$  with respect to  $\theta$  is equivalent to minimizing  $\sum_i (\sum_j y_{ij}^2 - m_i \bar{y}_i^2) + \sum_i m_i (\bar{y}_i - f_\theta(x_i))^2$  with respect to  $\theta$ . Since the first term does not depend on  $\theta$ , this is further equivalent to minimizing the weighted least squares RSS  $\sum_i m_i (\bar{y}_i - f_\theta(x_i))^2$  with respect to  $\theta$ .

7. (a) For linear regression, the estimator  $\hat{f}(x_0)$  has least squares estimate  $x_0^T (X^T X)^{-1} X^T Y$ .  $x_0^T (X^T X)^{-1} X^T$  is a  $1 \times N$  row vector whose entries depend only on  $x_0$  and  $\mathcal{X}$ , so the inner product  $(x_0^T (X^T X)^{-1} X^T) Y$  satisfies the required format.  
For kNN,  $\ell_i$  is  $1/k$  if  $x_i$  is one of the closest  $k$  neighbors to  $x_0$  and 0 otherwise. Again, this function depends only on  $x_0$  and  $\mathcal{X}$  so it satisfies the required format.
- (b)  $E_{Y|X}((f(x_0) - \hat{f}(x_0))^2) = \text{Var}_{Y|X}(\hat{f}(x_0)) + E_{Y|X}(f(x_0) - \hat{f}(x_0))^2 = \text{Var}_{Y|X}(\hat{f}(x_0)) + (f(x_0) - E_{Y|X}(\hat{f}(x_0)))^2$  since  $f(x_0)$  does not depend on the training data so can be considered constant with respect to the expectations/variances.
- (c) Using exactly analogous reasoning as the previous part, we obtain  $\text{Var}_{Y,X}(\hat{f}(x_0)) + (f(x_0) - E_{Y,X}(\hat{f}(x_0)))^2$ .
- (d) By Adam's law,  $E_{Y,X}(-) = E_X(E_{Y|X}(-))$  so the (c) result is the expectation over  $X$  of the (b) result.

8.

9. Let  $\tilde{\beta}$  be the least squares estimate for the test data. Since the testing and training data are drawn from the same underlying distribution/population, the least squares estimates  $\hat{\beta}$  and  $\tilde{\beta}$  are identically distributed. (In particular, extra data points do not result in scaling of either estimate.) Thus all the terms  $S_i = (y_i - \beta^T x_i)^2$  and  $\tilde{S}_i = (\tilde{y}_i - \beta^T \tilde{x}_i)^2$  in the sums in  $R_{tr}$  and  $R_{te}$  are identically distributed (though not

independent) with expectation  $E(S_i) = E(\tilde{S}_i) = S$  for some  $S$ . Thus  $E(R_{tr}(\hat{\beta})) = E(1/N \sum_i S_i) = 1/N \sum_i E(S_i) = S = 1/M \sum_i E(\tilde{S}_i) = E(1/M \sum_i \tilde{S}_i) = E(R_{te}(\tilde{\beta}))$ .

Next, by definition of least squares,  $R_{te}(\tilde{\beta})$  is minimized at  $\tilde{\beta}$  and in particular  $R_{te}(\tilde{\beta}) \leq R_{te}(\hat{\beta})$ . By monotonicity of expectation,  $E(R_{te}(\tilde{\beta})) \leq E(R_{te}(\hat{\beta}))$ . Thus  $E(R_{tr}(\hat{\beta})) = E(R_{te}(\tilde{\beta})) \leq E(R_{te}(\hat{\beta}))$