

Mock Midterm Exam - Nov 19, 2018

1 Subgradient Descent

Derive the (sub)gradient descent update rule for a one-parameter linear model using the Mean Absolute Error,

$$\mathcal{L}_{\text{MAE}}(\mathbf{X}, \mathbf{y}, w) = \frac{1}{N} \sum_{n=1}^N |wx_n - y_n|.$$

Hint: The function $f(x) = |ax|$ is a composition of two simpler functions. Use the chain rule!

The function can be seen as follows

$$f(v) = |v|$$

$$f'(v) = |v| \text{ sgn } v$$

using the chain rule

$$\text{we have that } v = wx_n - y_n$$

$$\frac{dv}{dw} = x_n$$

now

$$|v|' \left\{ \begin{array}{l} \text{if } v > 0 \\ |v|' = +1 \\ \text{if } v = 0 \\ |v|' = \text{subgradient } \approx 0 \\ \text{if } v < 0 \\ |v|' = -1 \end{array} \right.$$

so it can be seen as the sign of the function

$$f'(w) = \underline{\text{sgn } x_n}$$

2 Multiple-Output Regression

Let $S = \{(\mathbf{y}_n, \mathbf{x}_n)\}_{n=1}^N$ be our training set for a regression problem with $\mathbf{x}_n \in \mathbb{R}^D$ as usual. But now $\mathbf{y}_n \in \mathbb{R}^K$, i.e., we have K outputs for each input. We want to fit a linear model for each of the K outputs, i.e., we now have K regressors $f_k(\cdot)$ of the form

$$f_k(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}_k,$$

where each $\mathbf{w}_k^\top = (w_{k1}, \dots, w_{kD})$ is the weight vector corresponding to the k -th regressor. Let \mathbf{W} be the $D \times K$ matrix whose columns are the vectors \mathbf{w}_k .

Our goal is to minimize the following cost function \mathcal{L} :

$$\mathcal{L}(\mathbf{W}) = \sum_{k=1}^K \sum_{n=1}^N \frac{1}{2\sigma_k^2} (y_{nk} - \mathbf{x}_n^\top \mathbf{w}_k)^2 + \frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k\|_2^2,$$

where the σ_k are known real-valued scalars. Let $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_K)$.

For the solution, let \mathbf{X} be the $N \times D$ matrix whose rows are the feature vectors \mathbf{x}_n .

1. Write down the normal equations for \mathbf{W}^* , the minimizer of the cost function. I.e., what is the first-order condition that \mathbf{W}^* has to fulfill in order to minimize $\mathcal{L}(\mathbf{W})$.
2. Is the minimum \mathbf{W}^* unique? Assuming it is, write down an expression for this unique solution.
3. Write down a probabilistic model, so that the MAP solution for this model coincides with minimizing the above cost function. Note that this will involve specifying the likelihoods as well as a suitable prior (which will give you the regression term).

$$\begin{aligned} & \sum_{k=1}^K \left(\frac{1}{2\sigma_k^2} (y_{1k} - \mathbf{x}^\top \mathbf{w}_k)^2 + \frac{1}{2} \|\mathbf{w}_k\|_2^2 \right) \\ & - \frac{1}{2\sigma_k^2} \mathbf{x}^\top (y_{1k} - \mathbf{x}^\top \mathbf{w}_k) + \mathbf{w}_k = 0 \\ & \frac{\mathbf{x}^\top}{2\sigma_k^2} (\mathbf{x}^\top \mathbf{w}_k - y_{1k}) + \mathbf{w}_k = 0 \end{aligned}$$

$$\begin{aligned} \mathbf{x}^\top \mathbf{w}_k - \mathbf{x}^\top y_{1k} + \sigma_k^2 \mathbf{w}_k &= 0 \\ (\mathbf{x}^\top \mathbf{x} + \sigma_k^2) \mathbf{w}_k &= \mathbf{x}^\top y_{1k} \\ \mathbf{w}_k &= (\mathbf{x}^\top \mathbf{x} + \sigma_k^2)^{-1} \mathbf{x}^\top y_{1k} \end{aligned}$$

*The function is STRICTLY CONVEX
 \therefore It's a unique minimizer.

$$\begin{aligned} P(y_k | \mathbf{x}, \mathbf{w}_k) &= \mathbf{w}_k^\top \mathbf{x} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma_k^2) \\ \arg \min_{\mathbf{w}} \log \left(\prod_{k=1}^K \frac{1}{2\sigma_k^2} e^{-\frac{(y_k - \mathbf{w}_k^\top \mathbf{x})^2}{2\sigma_k^2}} \right) \end{aligned}$$

$$\mathcal{L}(\mathbf{w}) = \sum \frac{1}{2\sigma_k^2} \|y_k - \mathbf{w}_k^\top \mathbf{x}\|^2$$

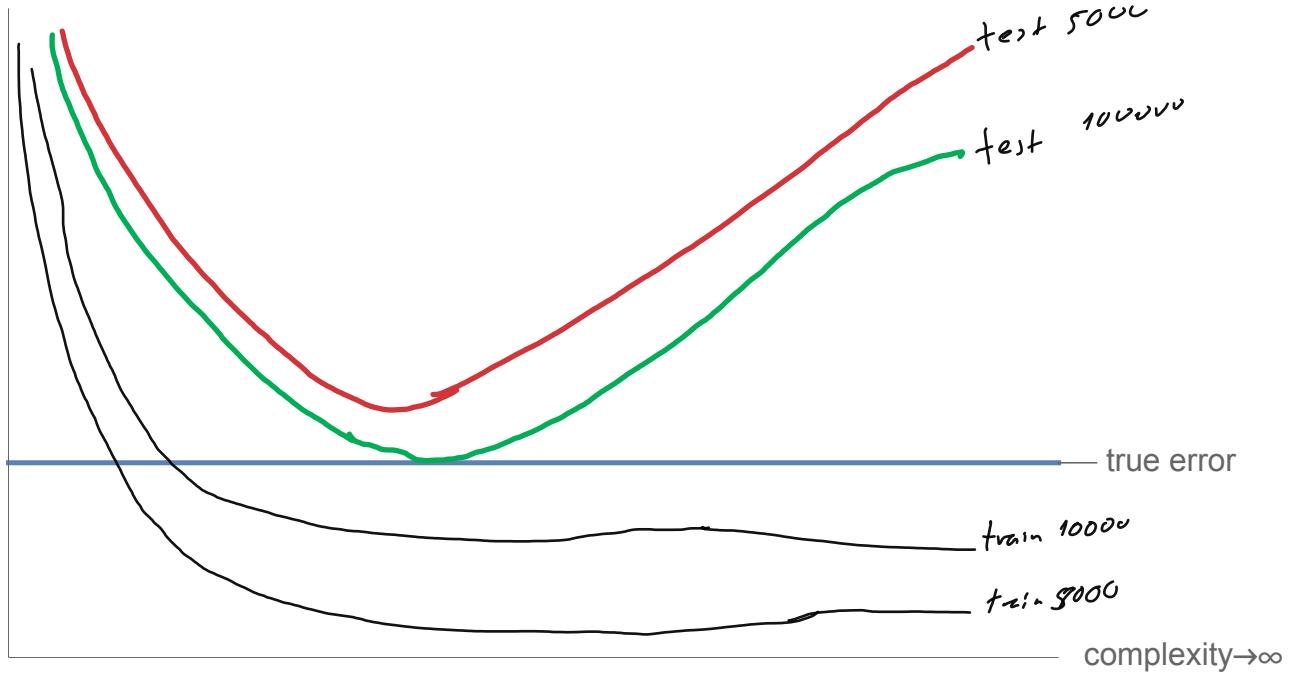
$$\begin{aligned} P(y | \mathbf{x}, \mathbf{w}) &= \prod_{k=1}^K \mathcal{N}(y_k | \mathbf{w}_k^\top \mathbf{x}, \sigma_k^2) \\ P(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | 0, \mathbf{I}) \end{aligned}$$

3 Bias Variance Trade-off (Due to Alex Smola)

Assume that you have two data sets that contain iid samples from the same distribution, call them S_1 and S_2 . S_1 contains 5000 samples, whereas S_2 contains 100000 samples. You randomly split each of the data sets into a training and a testing set, where eighty percent of the data is assigned to the training set. You then train and test on a family of increasing complexity.

In the figure below draw four curves, two that show the *training error* as a function of the model complexity (for S_1 and S_2) and two that show the *testing error* as a function of the model complexity (for S_1 and S_2). Label each of the 4 curves clearly. The constant curve labeled "true error" corresponds to the error due to the inherent noise in the samples and is drawn as a reference curve.

error $\uparrow \infty$



4 Exponential Families

Consider the Poisson distribution with parameter λ . It has a probability mass function given by $p(i) = \frac{\lambda^i e^{-\lambda}}{i!}$, $i = 0, 1, \dots$.

(i) Write $p(i)$ in the form of an exponential distribution $p(i) = h(i)e^{\eta\phi(i)-A(\eta)}$. Explicitly specify h , η , ϕ , and $A(\eta)$.

(ii) Compute $\frac{dA(\eta)}{d\eta}$ and $\frac{d^2A(\eta)}{d\eta^2}$? Is this the result you expected?

$$\begin{aligned} & \exp(\ln(\lambda^i) - \ln(i!) - \lambda) \\ & \exp\{i\ln(\lambda) - \lambda - \ln(i!)\} \\ & A(\eta) = e^\eta \quad h(i) = \frac{1}{i!} \\ & \eta = \ln(\lambda) ; \phi(i) = i \\ & \lambda = e^\eta \\ & \frac{dA(\eta)}{d\eta} = e^\eta = \lambda = \mu \\ & \frac{d^2A(\eta)}{d\eta^2} = e^\eta = \lambda = \sigma^2 \end{aligned}$$

5 Multiple Choice Questions and Simple Problems

Mark the correct **answer(s)**. More than one answer can be correct!

- In regression, “complex” models tend to

1. overfit
2. have large bias
3. have large variance

- In regression, “simple” models tend to

1. overfit
2. have large bias
3. have large variance

- We are given a data set $S = \{(\mathbf{x}_n, y_n)\}$ for a binary classification task where \mathbf{x}_n in \mathbb{R}^D . We want to use a *nearest-neighbor* classifier. In which of the following situations do we have a reasonable chance of success with this approach? [Ignore the issue of complexity.]

1. n is fixed, $D \rightarrow \infty$

2. $n = D^2$, $D \rightarrow \infty$

3. $n \rightarrow \infty$, $D \ll \ln(n)$

4. $n \rightarrow \infty$, D is fixed

- We add a regularization term because

1. this sometimes renders the minimization problem of the cost function into a strictly convex/concave problem

2. this tends to avoid overfitting

3. this converts a regression problem into a classification problem

- The k -nearest neighbor classifier

1. typically works the better the larger the dimension of the feature space

2. can classify up to k classes

3. typically works the worse the larger the dimension of the feature space

4. can only be applied if the data can be linearly separated

5. has a misclassification rate of at most two times the one of the Bayes classifier if we have lots of data

6. has a misclassification rate that is two times better than the one of the Bayes classifier

- A real-valued scalar Gaussian distribution

1. is a member of the exponential family with one scalar parameter

2. is a member of the exponential family with two scalar parameters

3. is not a member of the exponential family

- Which of the following statements is correct, where we assume that all the stated minima and maxima are in fact taken on in the domain of relevance.

1. $\max\{0, x\} = \max_{\alpha \in [0,1]} \alpha x$

2. $\min\{0, x\} = \min_{\alpha \in [0,1]} \alpha x$

3. Let $g(x) := \min_y f(x, y)$. Then $g(x) \leq f(x, y)$

4. $\max_x g(x) \leq \max_x f(x, y)$

5. $\max_x \min_y f(x, y) \leq \min_y \max_x f(x, y)$

- Which of the following statements are correct?

1. The training error is typically smaller than the test error.

2. The SVM (support vector machine) formulation we discussed can be optimized using SGD.

3. One iteration of SGD for ridge regression costs roughly $\Theta(ND)$, where N is the number of samples and D is the dimension.

4. Logistic regression as formulated in class can be optimized using SGD.

- You have given the 2D data shown in Figure 1. You are allowed to add one component to your data (in addition to a constant component) and then must use a linear classifier. What component should you pick?

1. $x_1 + x_2$

2. $1/|x_1 + x_2|$

3. $x_1 + 4x_2$

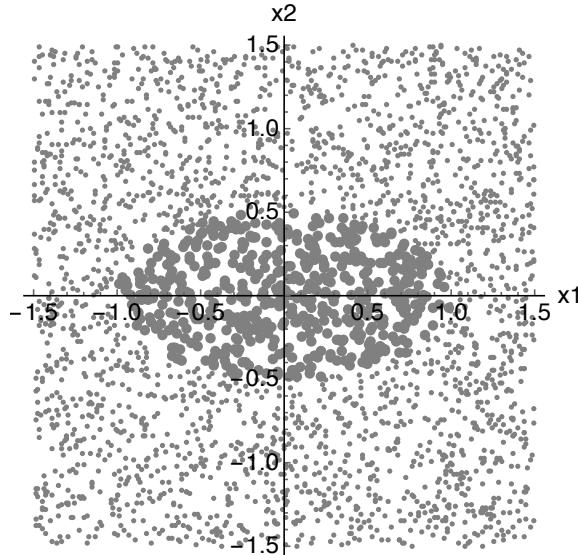


Figure 1: Some 2D data for classification. The two classes are indicated by different point sizes.

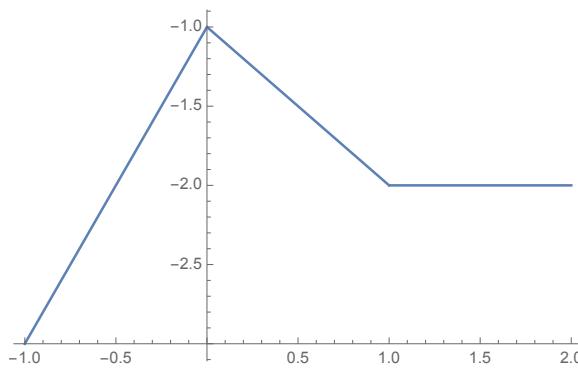


Figure 2: What is the subgradient of this function at $x = 1$?

4. $4x_1 + x_2$
5. $4x_1^2 + x_2^2$
6. $x_1^2 x_2^2$
7. $x_1^2 + 4x_2^2$

- The following functions are convex:

1. $f(x) := x^2, x \in \mathbb{R}$
2. $f(x) := x^3, x \in [-1, 1]$
3. $f(x) := -x^3, x \in [-1, 0]$
4. $f(x) := e^{-x}, x \in \mathbb{R}$
5. $f(x) := e^{-x^2/2}, x \in \mathbb{R}$
6. $f(x) := \ln(1/x), x \in (0, \infty)$
7. $f(x) := g(h(x)), x \in \mathbb{R}$, where g, h are convex and increasing over \mathbb{R}

- Let $f : \mathbb{R}^D \rightarrow \mathbb{R}$ be the function $f(\mathbf{w}) := \exp(\mathbf{x}^\top \mathbf{w})$, where $\mathbf{x} \in \mathbb{R}^D$. What is $\nabla_{\mathbf{w}} f$?
- Which of the following scalars g is a subgradient for the function shown in Figure 2 at the point $x = 1$?
 1. $g = -1$
 2. $g = -\frac{1}{2}$
 3. none exists
 4. $g = 0$