

1. **Multiclass classification** Consider the multiclass logistic regression optimization problem

$$\underset{\Theta \in \mathbb{R}^{n \times K}}{\text{maximize}} \quad f(\Theta) = \frac{1}{m} \sum_{i=1}^m \left(\sum_{k=1}^K y_{ik} x_i^T \theta_k - \log \sum_{k=1}^K \exp(x_i^T \theta_k) \right).$$

where $y_{ik} = 1$ if data sample i is in class k , and 0 otherwise. As usual, $x_i \in \mathbb{R}^n$ is the i th data feature. Here, we write the entire matrix variable as

$$\Theta = [\theta_1 \quad \theta_2 \quad \cdots \quad \theta_K].$$

(a) In terms of each θ_k , write the gradient of f with respect to θ_k . **Ans. (0.5 pts)**

$$\nabla_{\theta_k} f(\Theta) = \frac{1}{m} \sum_{i=1}^m \left(y_{ik} - \frac{\exp(x_i^T \theta_k)}{\sum_{l=1}^K \exp(x_i^T \theta_l)} \right) x_i, \quad \nabla_{\Theta} f(\Theta) = [\nabla_{\theta_1} f(\Theta) \quad \cdots \quad \nabla_{\theta_K} f(\Theta)].$$

(b) (Will be counted as extra credit.) Argue that this function has a smoothness parameter $L = L_X$ where $L_X = \frac{1}{m} \sum_{i=1}^m \|x_i\|_2^2$. Do so with the following substeps

- First, consider the function

$$g(s) = \log \left(\sum_{k=1}^K \exp(s_k) \right).$$

and find the Hessian of $g(s)$.

- Then, find the L -smoothness of g by showing that the maximum eigenvalue of $\nabla_s^2 g(s) \leq 1$.
- Then, for a sample point $x = x_i$, show that the maximum eigenvalue of $\nabla_{\Theta}^2 \hat{f}(\Theta) \leq \|x\|_2^2$ where

$$\hat{f}(\Theta) = \log \left(\sum_{l=1}^K \exp(x^T \theta_l) \right) = g(\Theta^T x).$$

- Use Jensen's inequality to reach the desired conclusion.

Ans. (1 pts) First, consider the function

$$g(s) = \log \left(\sum_{k=1}^K \exp(s_k) \right).$$

Then

$$\begin{aligned} \nabla_s g(s) &= \frac{1}{\sum_{l=1}^K \exp(s_l)} \exp(s) \\ \nabla_s^2 g(s) &= \frac{1}{\sum_{l=1}^K \exp(s_l)} \mathbf{diag}(\exp(s)) - \frac{1}{\left(\sum_{l=1}^K \exp(s_l) \right)^2} \exp(s) \exp(s)^T \end{aligned}$$

where for a vector $s \in \mathbb{R}^K$, $\exp(s) = [\exp(s_1), \dots, \exp(s_K)]^T$. We can find the largest eigenvalue of $\nabla_s^2 g(s)$ using the vector induced norm trick

$$\begin{aligned} \lambda_{\max}(\nabla_s^2 g(s)) &= \max_{\|u\|_2=1} u^T \nabla_s^2 g(s) u \\ &= \frac{1}{\sum_{l=1}^K \exp(s_l)} \underbrace{\max_{\|u\|_2=1} \sum_{i=1}^K u_i^2 \exp(s_i)^2}_{=\max_i \exp(s_i)^2} - \underbrace{\frac{(\exp(s)^T u)^2}{\sum_{l=1}^K \exp(s_l)}}_{\geq 0} \\ &\leq \frac{\max_i \exp(s_i)}{\sum_{l=1}^K \exp(s_l)} \\ &\leq 1. \end{aligned}$$

Then by chain rule, we can extrapolate further, by defining for a single x and taking $s = \Theta^T x$,

$$\begin{aligned}
\hat{f}(\Theta) &:= \log \left(\sum_{l=1}^K \exp(x^T \theta_l) \right) = g(\Theta^T x) \\
\nabla_{\theta_i \theta_j}^2 \hat{f}(\Theta) &= (\nabla_s^2 g(s))_{ij} x x^T \\
\lambda_{\max} \left(\nabla_{\Theta}^2 \hat{f}(\Theta) \right) &= \max_v \left\{ \sum_{i=1}^K \sum_{j=1}^K v_i^T \nabla_{\theta_i \theta_j}^2 \hat{f}(\Theta) v_j \mid \sum_{i=1}^K \|v_i\|_2^2 = 1 \right\} \\
&= \max_v \left\{ \sum_{i=1}^K \sum_{j=1}^K v_i^T x v_j^T x (\nabla_s^2 g(s))_{ij} \mid \sum_{i=1}^K \|v_i\|_2^2 = 1 \right\} \\
&= \max_{v, w} \left\{ w^T \nabla_s^2 g(s) w \mid \sum_{i=1}^K \|v_i\|_2^2 = 1, w_i = x^T v_i \right\} \\
&= \max_{v, w} \{ w^T \nabla_s^2 g(s) w \mid \|w\|_2 \leq \|x\|_2 \} \\
&\leq \|x\|_2^2
\end{aligned}$$

Therefore, for $f_i(\Theta) = \log \left(\sum_{l=1}^K \exp(x^T \theta_l) \right) = g(\Theta^T x_i)$,

$$\lambda_{\max} \left(\nabla_{\Theta}^2 \frac{1}{m} \sum_{i=1}^m f_i(\Theta) \right) \stackrel{\text{Jensen's ineq.}}{\leq} \lambda_{\max} \left(\nabla_{\Theta}^2 \hat{f}(\Theta) \right) \leq L_X$$

(c) The function

$$f(\theta) = \log \left(\sum_{i=1}^m \exp(\theta_i) \right)$$

is sometimes called the log-sum-exp function. As we saw in lecture, it has the nice property of acting like a soft-max function, by “pulling” away the largest values of θ_i , to somewhat exaggerate their “lead”.

A downside of using the log-sum-exp function is that it can have numerical issues. If θ_i is somewhat big, then $\exp(\theta_i)$ becomes very big, and can cause overflow. Conversely if θ_i is very negative, then all the values may be too close to 0 and cause underflow.

The “log-sum-exp-trick” is a numerical trick which deals with this issue, by adding and subtracting a constant whenever necessary. In effect, we simply do

$$f(\theta) = \log \left(\underbrace{\sum_{i=1}^m \exp(\theta_i - D)}_{f_1(\theta)} \right) + D.$$

Then, for the right choice of D , we can prevent overflow and underflow.

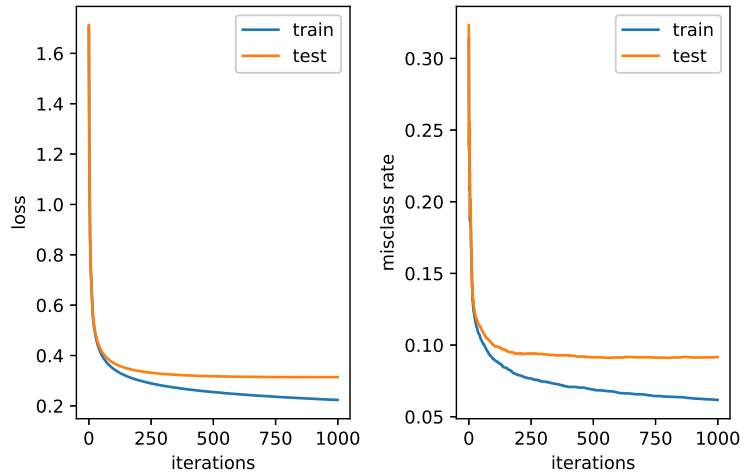
Propose a value of D such that $f_1(\theta) \leq 1$ (preventing overflow), and another value such that $f_1(\theta) \geq 1$ (preventing underflow).

Ans. (0.5 pts) Picking $D = \max_i \theta_i$ prevents overflow, and picking $D = \min_i \theta_i$ prevents underflow.

(d) **Coding.** Run multiclass logistic regression on MNIST dataset, against each of the 10 classes.

While usually we pick a stepsize of $2/L$, I have tried this and found a larger stepsize of 10^{-5} will work well. Use this stepsize and run for 500 iterations, or however many you need to see reasonable “working” behavior. Show the train/test loss plot and the train/test misclassification plot.

Ans. (1.0 pts)



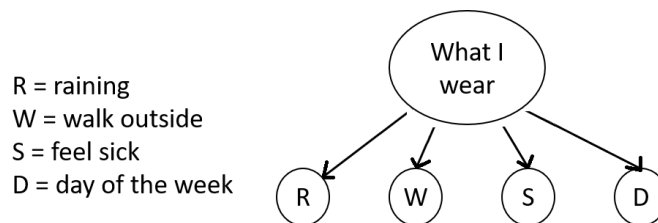
2. **Directed graphical models and probability inference** I have 4 tops: a red sweater, a blue T-shirt, a green hoodie, and a white tank top. I need your help to decide what to wear.

(a) I decide what to wear based on four factors:

- if it's raining
- if I want to take a walk outside
- if I feel sick
- the day of the week it is

Using the Naive Bayes assumption, draw a graphical model that indicates how I will make a decision of what to wear each day.

Ans. (1 pts)



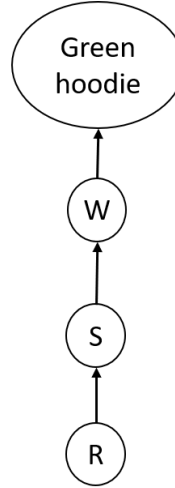
(b) To be more specific,

- I only wear the green hoodie when I walk outside, regardless of all other factors.
- If I feel sick, I will walk outside 10% of the time. If I feel well, I will walk outside 60% of the time.
- When it rains, I feel sick 70% of the time; otherwise, I feel sick 15% of the time.

Draw a corresponding graphical model for determining whether I wear a green hoodie. Given that it is raining, infer the probability that I am wearing a green hoodie.

Ans. (1 pts)

R = raining
W = walk outside
S = feel sick
D = day of the week



$$\begin{aligned}
 \Pr(\text{sick} = S | \text{rain} = 1) &= \begin{cases} 0.7 & \text{if } S = 1 \\ 0.3 & \text{if } S = 0 \end{cases} \\
 \Pr(\text{walk} = W | \text{rain} = 1) &= \Pr(\text{walk} = W | \text{sick} = 1) \Pr(\text{sick} = 1 | \text{rain} = 1) \\
 &\quad + \Pr(\text{walk} = W | \text{sick} = 0) \Pr(\text{sick} = 0 | \text{rain} = 1) \\
 &= \begin{cases} 0.1 \cdot 0.7 + 0.6 \cdot 0.3 = 0.25 & \text{if } W = 1 \\ 0.9 \cdot 0.7 + 0.4 \cdot 0.3 = 0.75 & \text{if } W = 0 \end{cases} \\
 \Pr(\text{green hoodie} = G | \text{rain} = 1) &= \Pr(\text{green hoodie} = G | \text{walk} = 1) \Pr(\text{walk} = 1 | \text{rain} = 1) \\
 &\quad + \Pr(\text{green hoodie} = G | \text{walk} = 0) \Pr(\text{walk} = 0 | \text{rain} = 1) \\
 &= \begin{cases} 1 \cdot 0.25 + 0 \cdot 0.75 & \text{if } G = 0 \\ 0 \cdot 0.25 + 1 \cdot 0.75 & \text{if } G = 1 \end{cases}
 \end{aligned}$$

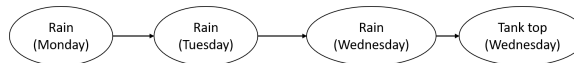
So, the probability that I will wear a green hoodie is 25%.

(c) The probability that I wear a tank top, independently of all the other clothes, is 75% if it's raining and 25% if it's not raining.

- Today is Monday and it is raining.
- The probability that it will rain, given that the previous day rained, is 70%. The probability that it will rain, given that the previous day did not rain, is 10%.

Draw a graphical model predicting whether I will wear a tank top on Wednesday, and calculate this probability.

Ans. (1 pts)



Denote M, T, W the random variables for the event of raining on Monday, Tuesday, and Wednesday. Denote U for the event that I wear the white tank top.

$$\begin{aligned}
\Pr(T|M=1) &= \begin{cases} 0.7 & \text{if } T=1 \\ 0.3 & \text{if } T=0 \end{cases} \\
\Pr(W|M=1) &= \Pr(W|T=1)\Pr(T=1|M=1) \\
&\quad + \Pr(W|T=0)\Pr(T=0|M=1) \\
&= \begin{cases} 0.7 \cdot 0.7 + 0.1 \cdot 0.3 = 0.52 & \text{if } W=1 \\ 0.3 \cdot 0.7 + 0.9 \cdot 0.3 = 0.48 & \text{if } W=0 \end{cases} \\
\Pr(U|M=1) &= \Pr(U|W=1)\Pr(W=1|M=1) + \Pr(U|W=0)\Pr(W=0|M=1) \\
&= \begin{cases} 0.75 \cdot 0.52 + 0.25 \cdot 0.48 = 0.51 & \text{if } U=0 \\ 0.25 \cdot 0.52 + 0.75 \cdot 0.48 = 0.49 & \text{if } U=1 \end{cases}
\end{aligned}$$

So, the probability that I will wear a white tank top is 51%.

3. **Hidden Markov Model spellchecker** In this exercise we will make a spell-checker using a HMM. To do this, download `alice_nlp_release.ipynb` and follow the instructions.

- Read through the first two blocks to get an idea of what the task is. The idea is to go through the corrupted corpus, identify words which have probably been corrupted, and correct them probabilistically.
- In the 4th box, fill in the functions to construct the word probabilities (weighted frequencies in uncorrupted corpus) and transition matrix (which gives $\Pr(\text{word} \mid \text{prev word})$). If done correctly, the lines printed out should read

```

prob. of "alice" 0.014548615047424706
prob. of "queen" 0.002569625514869818
prob. of "chapter" 0.0009069266523069947

with smoothing

prob. of "the alice" 0.00025406504065040653
prob. of "the queen" 0.016514227642276422
prob. of "the chapter" 0.012957317073170731

no smoothing

prob. of "the alice" 0.0
prob. of "the queen" 0.03968253968253968
prob. of "the chapter" 0.0
prob. of "the hatter" 0.031135531135531136

```

- In the 5th box, fill in the function for computing the emission probability. The first 10 words closest to Alice should be

```
['abide', 'alice', 'above', 'voice', 'alive', 'twice', 'thick', 'dance', 'stick', 'prize']
```

- Construct and run your Hidden Markov Model spell checker using the functions computed for the prior probabilities, emission probabilities, and transition probabilities. List some words whose spelling was corrected correctly, and some examples where the spell-correcter did not work as expected. Report the recovery rate of the “fixed” corpus.

Ans. (4 pts) Recovery rate of corrupted corpus: 0.759.

Recovery rate of fixed corpus, with smoothed transition matrix: 0.187. (Yikes! There was a bug in my code originally. This is not what I had hoped the recovery rate would be, but it turned out that smoothing was a terrible idea, since most of the values in the transition matrix should have been 0.)

Recovery rate of fixed corpus without smoothing: 0.924.

(Note: this is the only part of the problem that needs to be scored.)

Challenge!

How to deal with multiple advisers I have m advisers, and I don't know which one of them to trust. Every day $t = 1, \dots, T$, I ask all m advisers a prediction question, which they answer yes $y_i^{(t)} = 1$ or no $y_i^{(t)} = -1$ ($i = 1, \dots, m$). The true answer on each day is denoted by $y_\star^{(t)} \in \{-1, 1\}$.

Each day, I have to make a guess as to what $y_\star^{(t)}$ has to be, which I do using the linear predictor

$$\hat{y}_{\text{pred}}(w, y^{(t)}) = \text{sign} \left(\sum_{i=1}^m w_i y_i^{(t)} \right).$$

The variables w_i , $i = 1, \dots, m$ are positive weights I place on each adviser, to indicate my trust in that adviser. I constrain $0 \leq w_i \leq 1$ and require $\sum_{i=1}^m w_i = 1$ (w represents a probability mass function.)

To learn how much I trust each adviser, I optimize the following convex optimization problem, over the weights w_i , over T days of observation:

$$\begin{aligned} \underset{w}{\text{minimize}} \quad & f(w) = \sum_{t=1}^T \underbrace{\left(\sum_{i: y_i^{(t)} \neq y_\star^{(t)}} w_i - \sum_{i: y_i^{(t)} = y_\star^{(t)}} w_i \right)}_{=: f^{(t)}(w)} \\ \text{subject to} \quad & 0 \leq w_i \leq 1 \\ & \sum_{i=1}^m w_i = 1 \end{aligned} \tag{1}$$

1. **Online learning method.** To learn the correct weighting, I run the following method. I start by setting $w_i^{(0)} = 1/m$ for all i . Then, every day, if an adviser is correct, I update that adviser's weight as $\hat{w}_i^{(t)} = w_i^{(t-1)} e^\mu$ for some $\mu \geq 0$. If the adviser is not correct, I leave the weight alone $\hat{w}_i^{(t)} = w_i^{(t-1)}$. Then I reweight, e.g.

$$w_i^{(t)} = \frac{\hat{w}_i^{(t)}}{\sum_{j=1}^m \hat{w}_j^{(t)}}.$$

This method is known as the *exponential weighting method*.

Suppose that each adviser has a probability p_i of being correct. Show that as $T \rightarrow +\infty$, this method converges to exclusively listening to the adviser which is most often correct, e.g.

$$w_i^{(t)} \xrightarrow{t \rightarrow +\infty} \begin{cases} 1 & \text{if } i = \underset{j}{\operatorname{argmax}} p_j \\ 0 & \text{else.} \end{cases}$$

Hint: Note that the method doesn't really change if you reweight at each iteration, vs reweighting at the very end.

Ans. In this entire problem, it is often easier to consider a change of variables $y^{(t)} = \frac{1}{t} \log(w^{(t)})$ (log taken element-wise). Denoting $z_i^{(t)}$ the variable indicating each adviser's choice:

$$z_i^{(t)} = \begin{cases} 1 & \text{if } i\text{th adviser says yes on day } t \\ 0 & \text{else.} \end{cases}$$

Then $\mathbb{E}[z_i^{(t)}] = p_i$. Then

$$w_i^{(t)} = \exp \left(\mu \sum_{\tau=1}^t z_i^{(\tau)} \right) \iff y_i^{(t)} = \frac{\mu}{t} \sum_{\tau=1}^t z_i^{(\tau)}.$$

Here, we can see that $y_i^{(t)}$ is a random variable with mean μp_i and variance

$$\frac{\mu^2}{t} p_i (1 - p_i) \xrightarrow{t \rightarrow +\infty} 0.$$

So, $y_i^{(t)} \xrightarrow{t \rightarrow +\infty} \mu p_i$.

Now we follow the hint, and do not normalize after each step. Then

$$w_i^{(t)} = t \exp(y_i^{(t)}) \xrightarrow{t \rightarrow +\infty} \exp(\mu p_i).$$

Take $p_{\max} = \max_i p_i$, and scale the reweighting factor as

$$\tilde{C} = C e^{t p_{\max} \mu}.$$

Then

$$w_i^{(t)} = \tilde{C} e^{t(p_i - p_{\max})\mu}.$$

If $p_i = p_{\max}$, then $w_i^{(t)} \rightarrow \tilde{C}$. Otherwise, if $p_i < p_{\max}$, then as $t \rightarrow +\infty$, $e^{t(p_i - p_{\max})\mu} \rightarrow 0$ and $w_i^{(t)} \rightarrow 0$.

2. **Projected incremental gradient descent.** We can also think of a gradient descent method to solve (??), where at each iteration, we take a gradient step

$$w^{(t+1)} = \mathbf{proj}_{\Delta_{m-1}} \left(w^{(t)} - \eta \nabla f^{(t)}(w^{(t)}) \right)$$

Show that taking a gradient step

$$\hat{w} = w^{(t)} - \eta \nabla f^{(t)}(w^{(t)})$$

essentially results in subtracting or adding a constant value to $w^{(t)}$ each time an adviser gets an answer wrong or right, respectively.

Ans. Basically, the point is to derive

$$\nabla f^{(t)}(w^{(t)})_i = -y_*^{(t)} y_i^{(t)}$$

and thus

$$-\eta \nabla f^{(t)}(w^{(t)})_i = \begin{cases} \eta & \text{if } y_i^{(t)} = y_*^{(t)} \\ -\eta & \text{else.} \end{cases}$$

3. **Mirror descent.** The set $\Delta_{m-1} = \{w : 0 \leq w \leq 1, \sum_{i=1}^m w_i = 1\}$ is called the *probability simplex*. Projecting on this set is actually not very trivial, and involves sorting all the elements (which, if we remember from algorithms, is at least $O(m \log m)$). So, while we could run a projected gradient descent method to solve (??), we will try a different trick instead.

We do something called the *mirror descent method*. Basically, this method is a projected gradient method, but with a transformed variable. In particular, our *mirror map* is going to be the gradient of the strictly convex function

$$g(u) = \sum_{i=1}^m u_i \log(u_i) \quad (\text{negative entropy function}).$$

Derive the mirror map, which is $\nabla g(u)$, and the inverse mirror map. That is, find Φ and Φ^{-1} where

$$\Phi(u) = \nabla g(u), \quad \Phi^{-1}(\nabla g(u)) = u.$$

Ans.

$$\Phi(u) = \begin{bmatrix} 1 + \log(u_1) \\ \vdots \\ 1 + \log(u_m) \end{bmatrix}, \quad \Phi^{-1}(v) = \begin{bmatrix} e^{v_1 - 1} \\ \vdots \\ e^{v_m - 1} \end{bmatrix},$$

4. The incremental mirror descent method can then be summarized as

$$\begin{aligned} \Phi(\hat{w}) &= \Phi(w^{(t)}) - \nabla f(w^{(t)}) \\ \Phi(w^{(t+1)}) &= \Phi(\mathbf{proj}_{\Delta_{m-1}}(\hat{w})). \end{aligned}$$

Or, to write it in a way that is more implementable,

$$\begin{aligned}\hat{w} &= \Phi^{-1}\left(\Phi(w^{(t)}) - \nabla f(w^{(t)})\right) \\ w^{(t+1)} &= \Phi^{-1}\left(\Phi(\mathbf{proj}_{\Delta_{m-1}}(\hat{w}))\right).\end{aligned}$$

Show that this method is equivalent to the exponential weighted algorithm in part (a).

Ans. This problem is a mistake on my end, in that the projection on the simplex is in fact not equivalent to reweighting. Rather, it is a much more involved procedure that involves sorting. The projection I should have written is with respect to the KL divergence; that is, for $w = \frac{1}{\hat{w}^T \mathbf{1}} \hat{w}$,

$$w = \operatorname{argmin}_w \left\{ -\sum_i \hat{w}_i^T \log(w_i) \mid w^T \mathbf{1} = 1, w \geq 0 \right\} =: \mathbf{proj}_{\Delta}^{KL}(\hat{w}).$$

Anyway, I gave you guys points if you wrote something reasonable, ignoring the projection step. Apologies for confusion!

Another way to write the first step is as

$$\log(\hat{w}) = \log(w^{(t)}) - \eta y_*^{(t)} y^{(t)}$$

which can be rearranged to

$$\hat{w} = w^{(t)} e^{-\eta y_*^{(t)} y^{(t)}}$$

Since \hat{w} is always nonnegative, projecting on the unit simplex simply involves rescaling, so that

$$w^{(t+1)} = \mathbf{proj}_{\Delta_{m-1}}^{KL}(\hat{w}) = \frac{1}{\hat{w}^T \mathbf{1}} \hat{w}.$$

This is exactly the exponential weighting method!