

Statistical Methods: Portfolio 3

By Henry Bourne

1 Linear Classifiers

In the portfolio 1 we saw how to conduct binary classification, now we will look at how to conduct **multi-class classification**. This is where we have an input $\mathbf{x} \in \mathbb{R}^d$ and an output $y \in \{1, \dots, K\}$.

The geometry of the problem in this case is more complicated than we had with binary classification, we can no longer simply check the sign of a single $f(\mathbf{x})$ to classify. We could try introducing multiple functions. Let's say we have 3 classes, we could try introducing another function so now we have f_1 and f_2 , we could then perform classification by checking the signs of both f_1 and f_2 and have this dictate the class. However, this can get confusing as for example in this case we have 4 possible outcomes, $\{+, -\}^2$, for the signs of our f 's, however, we only have 3 classes. What about if we introduce pairwise binary classifiers, ie. we have $f_{i,j}$ classifies a point as either i or j and we have such a function for any pair of classes, then we classify a point as the majority vote given by the binary classifiers. The problem here is that all the classifiers may disagree in which case there is no majority vote.

Rather than relying on the sign of a function f to make predictions lets instead fit a vector valued function $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^K$, where K is the number of classes. Given an input \mathbf{x} our prediction is $\hat{k} = \operatorname{argmax}_k f^{(k)}(\mathbf{x})$. Note that this no longer has a simple geometric interpretation anymore.

1.1 Least Squares Classifier

We will first define the least squares classifier in the binary classification case and after extend it to the multi-class case:

Definition 1.1 Least Squares Binary Classifier:

We first perform LS on the data, ie. find:

$$\mathbf{w}_{LS} := \operatorname{argmin}_{\mathbf{w}} \sum_{i \in D} [y_i - f(\mathbf{x}_i; \mathbf{w})]^2 \quad (1)$$

Then we can find the predicted label $\hat{y} := \operatorname{sign}(f(\mathbf{x}_i; \mathbf{w}_{LS}))$

Note that we can also use a feature transform for f as well, which would allow us to fit more complex classifiers, as data not separable in the original space may be separable in the feature space. In the multi-class case:

Definition 1.2 Multi-class LS classification:

We use a **one-hot encoding** which is where we replace $y_i = k$ in our data with $\mathbf{t}_i \in \{0, 1\}^K$ where all entries are 0 bar $t_i^{(k)} = 1$. Then we have:

$$\mathbf{W}_{LS} := \operatorname{argmin}_{\mathbf{W}} \sum_{i \in D} \|\mathbf{t}_i - f(\mathbf{x}_i; \mathbf{W})\|^2 \quad (2)$$

where $\mathbf{W} \in \mathbb{R}^{(d+1) \times K}$, $f(\mathbf{x}; \mathbf{W}) = \mathbf{W}^T \tilde{\mathbf{x}}$ and $\tilde{\mathbf{x}} := [\mathbf{x}^T, 1]^T$.

Then our prediction $\hat{y} = \operatorname{argmax}_k f(\mathbf{x}; \mathbf{W})^{(k)} = \operatorname{argmax}_k (\mathbf{w}_{LS}^{(k)})^T \tilde{\mathbf{x}}$, where $\mathbf{w}^{(k)}$ is the k -th column of \mathbf{W} .

Although this method can work the square loss tends not to make sense in a classification task, as a point far away from the boundary (fit without that point) can dramatically affect our decision boundary, even if it is correctly classified by the decision boundary (fit without that point). Also, unlike with LS regression, LS classification lacks a probabilistic interpretation.

1.2 Fisher Discriminant Analysis (FDA)

Note that taking the inner product $\langle \mathbf{w}, \mathbf{x} \rangle$ embeds \mathbf{x} onto a one-dimensional line along the \mathbf{w} direction. We can say \mathbf{w} gives a good embedding if the \mathbf{x} it embeds are close together in the embedding if they are from the same class but far apart if they are from different classes. We can define these two properties we want from our embedding as:

Definition 1.3 Within-class Scatterness:

For embedding $\mathbf{w}^T \mathbf{x}$, the **embedding centre** for class k is:

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i, y_i=k} \mathbf{w}^T \mathbf{x}_i \quad (3)$$

then the within-class scatterness of class k is:

$$s_{\mathbf{w},k} := \sum_{i, y_i=k} (\mathbf{w}^T \mathbf{x}_i - \hat{\mu}_k)^2 \quad (4)$$

Definition 1.4 Between-class Scatterness:

The **embedded dataset centre** is:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i \quad (5)$$

then the between-class scatterness is:

$$s_{b,k} = n_k (\hat{\mu}_k - \hat{\mu})^2 \quad (6)$$

note the n_k is needed to make $s_{b,k}$ the same scale as $s_{\mathbf{w},k}$.

Then ideally we would like to maximize the between-class scatterness and minimize the within-class scatterness for all the classes, which is what the following does:

Definition 1.5 Fisher Discriminant Analysis (FDA):

$$\max_{\mathbf{w}} \left[\sum_k s_{b,k} / \sum_k s_{\mathbf{w},k} \right] \quad (7)$$

if $K = 2$ then this has a simple solution: $\mathbf{w} := \mathbf{S}_{\mathbf{w}}^{-1}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-)$, where $\mathbf{S}_{\mathbf{w}} := \sum_{k=1}^K \mathbf{S}_k$ and \mathbf{S}_k is the sample covariance matrix of class k times n_k . However, note that the FDA does not learn a decision function f , the \mathbf{w}_{FDA} obtained cannot be used directly by the prediction function to make a prediction. This is because (eg. in the binary case) $f(\mathbf{x}; \mathbf{w}_{FDA}) > 0$ does not mean that \mathbf{x} is predicted as the positive class. The FDA also doesn't care about classification accuracy, ie. minimizing the FP or FN rate.

1.3 Probabilistic (Generative) Classifiers

How do we put a classification problem under a probabilistic framework? Let's say we would like to **minimize the expected loss**: $\hat{y} := \operatorname{argmin}_{y_0} \mathbb{E}_{p(y|\mathbf{x})}[L(y, y_0)|\mathbf{x}]$. To do this we need $p(y|\mathbf{x})$. The **discriminative approach** would be to infer $p(y|\mathbf{x})$ directly. Here, we will look at the **generative approach** which is where we note that $p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$ and we infer $p(\mathbf{x}|y)$ ¹.

1.3.1 Continuous input

If \mathbf{x} is a continuous variable then the Multi-Variate Normal (MVN) is a natural choice for the model of $p(\mathbf{x}|y)$, ie. we have $p(\mathbf{x}|y=k; \mathbf{w}) := N_{\mathbf{x}}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ (assuming iid data and that all classes have same covariance matrix). We can write the likelihood of the parameters over the data, D , as: $p(D|\mathbf{w}) = \prod_{i \in D} p(\mathbf{x}_i, y_i|\mathbf{w}) = \prod_{i \in D} p(\mathbf{x}_i|y_i; \mathbf{w})p(y_i) = \prod_{i \in D} N_{\mathbf{x}_i}(\boldsymbol{\mu}_{y_i}; \boldsymbol{\Sigma})p(y_i)$. Therefore our MLE's are: $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \hat{\boldsymbol{\Sigma}} := \operatorname{argmax}_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \hat{\boldsymbol{\Sigma}}} \sum_{i \in D} \log[N_{\mathbf{x}_i}(\boldsymbol{\mu}_{y_i}; \boldsymbol{\Sigma})p(y_i)]$. If we plug in our estimates for $p(y_i = k)$ which are $\frac{n_k}{n}$, then we can find the MLE for $\hat{\boldsymbol{\mu}}_k := \frac{1}{n_k} \sum_{i \in D, y_i=k} \mathbf{x}_i$. We can then plug in our $\hat{\boldsymbol{\mu}}_k$'s to work out: $\hat{\boldsymbol{\Sigma}} := \sum_{k=1, \dots, K} \frac{n_k}{n} \frac{1}{n_k} \sum_{i \in D, y_i=k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T$,

¹ $p(y)$ is just the probability of the class being class y

where we notice that in red is the MLE of the covariance of individual classes. Our prediction is $\hat{y} := \operatorname{argmax}_y p(y|\mathbf{x}; \hat{\mathbf{w}}) \propto p(\mathbf{x}|y; \hat{\mathbf{w}})p(y)$. We can also prove that when using the shared covariance matrix MVN model, the decision boundary is piecewise-linear (Section A.1).

What if we assume that for each class k there are different covariance matrices? then the MLE reduces to estimating individual $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$. Note that in this case the decision boundary is no longer linear.

1.3.2 Discrete input

Here we will look at the case where \mathbf{x} is discrete. Let $\mathbf{x} := [x^{(1)}, \dots, x^{(d)}]^T$ and assume $x^{(1)}, \dots, x^{(d)}$ follow a multinomial distribution, where each $x^{(i)}$ is a quantity for a feature. Using Bayes rule we have that $p(y = k|\mathbf{x}) = (p(\mathbf{x}|y = k) \cdot p(y = k))/p(\mathbf{x})$. We now introduce our "naive assumptions" which are that all the features in \mathbf{x} are mutually independent, conditional on $y = k$. Using these assumptions we can write $p(y = k|\mathbf{x}) \propto p(y = k) \prod_{i=1}^n p(x_i|y = k)$ ². Therefore $p(\mathbf{x} = \mathbf{x}_0|y = k) \propto \prod_{i=1, \dots, d} \beta(i|y = k)^{x_0^{(i)}}$, for some \mathbf{x}_0 where $\beta(i|y = k)$ is the probability feature i occurs in class k . It is easy to estimate this quantity:

$$\beta(i|y = k) \approx \frac{\sum_{j \in D, y_j = k} x_j^{(i)}}{\sum_{j \in D, y_j = k} \sum_{l=1}^d x_j^{(l)}} \quad (8)$$

Then our prediction is $\hat{y} := \operatorname{argmax}_y p(\mathbf{x} = \mathbf{x}_0|y) \cdot p(y)$ where as before $p(y = k)$ can be estimated as $\frac{n_k}{n}$ and note that $\beta(i|y = k)$ can be estimated by counting. This is called **Naive Bayes classification**.

2 Discriminative Classifiers

Before we looked at Generative Classifiers, now we will look at discriminative classifiers. **Discriminative classifiers** aim to infer $p(y|\mathbf{x})$ given dataset D . The process normally looks like so:

1. Make a model assumption $p(y|\mathbf{x}; \mathbf{w})$.
2. Construct the likelihood function $p(D|\mathbf{w})$.
3. Estimate the parameters (eg. using MLE, MAP, full probabilistic approach, etc.).

The first question we ask is: what model should we use? Bayes rule says:

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|y)p(y)}{\sum_{y'} p(\mathbf{x}, y')} = \frac{p(\mathbf{x}|y)p(y)}{\sum_{y'} p(\mathbf{x}|y')p(y')} \quad (9)$$

Notice that we are now representing $p(y|\mathbf{x})$ using $p(\mathbf{x}|y)$. Suppose $y \in \{-1, 1\}$ then $p(y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|y=1)p(y=1)}{p(\mathbf{x}|y=1)p(y=1) + p(\mathbf{x}|y=-1)p(y=-1)}$. Assume $p(\mathbf{x}|y)p(y) > 0, \forall \mathbf{x}, y$ then this equals $1/(1 + \frac{p(\mathbf{x}|y=-1)p(y=-1)}{p(\mathbf{x}|y=1)p(y=1)})$, ie. we can rewrite $p(y|\mathbf{x})$ using the ratio $\frac{p(\mathbf{x}|y=-1)p(y=-1)}{p(\mathbf{x}|y=1)p(y=1)}$. So in a discriminative learning we model the **density ratio** $\frac{p(\mathbf{x}|y=-1)}{p(\mathbf{x}|y=1)}$ as opposed to the class density like in generative learning³. We will model the log of the inverse of the density ratio (as $f(\mathbf{x}; \mathbf{w})$) and so in our expression above we can replace the density ratio with $\exp(-f(\mathbf{x}; \mathbf{w}))$.

The **sigmoid function** is denoted $\sigma(t) := \frac{1}{1 + \exp(-t)}$. Then note that $p(y = 1|\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-f(\mathbf{x}; \mathbf{w}))} = \sigma(f(\mathbf{x}; \mathbf{w}))$, also note that when $y = -1$ our density ratio is inverse and so we no longer need the minus in front of f , therefore, in general $p(y|\mathbf{x}; \mathbf{w}) = \sigma(f(\mathbf{x}; \mathbf{w}) \cdot y)$.

Now, assuming the data, D , is iid., we would like to find the MLE for $p(y|\mathbf{x}; \mathbf{w})$:

$$\mathbf{w}_{MLE} = \operatorname{argmax}_{\mathbf{w}} \log \prod_{i \in D} p(y_i|\mathbf{x}_i; \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log p(y_i|\mathbf{x}_i; \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) \quad (10)$$

²We get this by discarding the $p(\mathbf{x})$, which makes it proportional and then using our assumption on $p(\mathbf{x}|y = k)$

³As a side note: clearly modelling the density ratio requires a lot fewer assumptions on your class densities, so models on class densities will model the density ratio, however models on the density ratio will not necessarily model the class density

This procedure is called **Logistic Regression** (even though it is not a regression!). Note that logistic regression is robust to outliers, unlike the LS classifier, this is because the sigmoid function tapers their affect. As before with the LS classifier we can also fit non-linear classifiers using a feature transform.

Apart from logistic regression there are other methods for estimating $p(y|\mathbf{x}; \mathbf{w})$. For example assuming a prior on \mathbf{w} we could use MAP:

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log(\sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) p(\mathbf{w})) = \operatorname{argmax}_{\mathbf{w}} \sum_{i \in D} \log(\sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i)) + \log(p(\mathbf{w})) \quad (11)$$

We can also use the full probabilistic approach ⁴:

$$p(y|\mathbf{x}) = \int p(y|\mathbf{x}; \mathbf{w}) p(\mathbf{w}|D) d\mathbf{w} \propto \int p(y|\mathbf{x}; \mathbf{w}) p(D|\mathbf{w}) p(\mathbf{w}) d\mathbf{w} \quad (12)$$

$$\propto \int \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) \prod_{i \in D} \sigma(f(\mathbf{x}_i; \mathbf{w}) \cdot y_i) p(\mathbf{w}) d\mathbf{w} \quad (13)$$

What about in the multi-class scenario? we can extend logistic regression to the multi-class scenario, where the density ratio is now:

$$p(y = i|\mathbf{x}) = \frac{p(\mathbf{x}|y = i)p(y = i)}{\sum_k p(\mathbf{x}|y = k)p(y = k)} \quad (14)$$

We expand further on this in conjunction with one-hot encoding in question B.2.3. Note that unlike in LS, logistic regression does not have a closed form solution and therefore must employ numerical methods to find \mathbf{w}_{MLE} .

3 Support Vector Machines

In binary classification it might be that there are many decision boundaries we could draw that would correctly classify all the training data, this leaves us asking: what is the "optimal" decision boundary in binary classification? For our decision boundary to be generalizable we want to minimize the error on unseen datasets rather than the training set. We would like a decision boundary that will correctly classify a small perturbation of one of our training data points whilst still correctly classifying all of our training data, ie. we want to maximize the gap (the **margin**) between the decision boundary and the lines given by $f(\mathbf{x}; \mathbf{w}) = 1$ and $f(\mathbf{x}; \mathbf{w}) = -1$, whilst keeping data points on the correct side of the margin.

If we are using decision boundary $f(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}_1, \mathbf{x} \rangle + w_0$ then we can calculate the thickness of our margin as $\frac{1}{\|\mathbf{w}_1\|}$. Then we can formulate the following constrained minimization problem: minimize $\|\mathbf{w}_1\|^2$ subject to $\forall i, y_i \cdot f(\mathbf{x}_i; \mathbf{w}) \geq 1$.

However, in many cases the dataset is not separable and so it is not possible to satisfy our constraint. We can relax these constraints which is what we do in the **soft-margin classifier**: $\min_{\mathbf{w}, \epsilon} \|\mathbf{w}_1\|^2 + \sum_i \epsilon^{(i)}$ subject to $\forall i, y_i(\langle \mathbf{w}_1, \mathbf{x}_i \rangle + w_0) + \epsilon^{(i)} \geq 1, \epsilon^{(i)} \geq 0$. It turns out that this classifier is actually a convex minimization problem and so every local minimum is a global minimum.

We can use the lagrangian dual (discussed in section A.2) to turn the minimization problem of the soft-margin classifier into: $\max_{\boldsymbol{\lambda}} -\frac{\tilde{\mathbf{X}}^T \mathbf{X}^T \mathbf{X} \tilde{\mathbf{X}}}{4} + \langle \boldsymbol{\lambda}, \mathbf{1} \rangle$ subject to $0 \leq \lambda_i \leq 1, \sum_i \lambda_i y_i = 0$. Note that the constraints here are simpler and that its quadratic wrt. $\boldsymbol{\lambda} \in \mathbb{R}^n$ as opposed to quadratic wrt. $\mathbf{w} \in \mathbb{R}^{d+1}$ in our previous optimization problem. Hence, the lagrangian dual is slow when n is large and our previous problem is slow when d is large, also note that it is possible to use kernel methods in solving the lagrangian!

Notice that the border of the (optimal) margins will always pass through some data points, these points can be thought of as resisting the expansion of the margin, hence are called **support vectors**. Hence the name of this method: "Support Vector Machines" (SVM). However, the SVM has some limitations, for example it is important to note the SVM is not a probabilistic classifier as it lacks a probabilistic interpretation. Also note that the computational cost of SVM is high as it involves solving a constrained optimization problem (as opposed to unconstrained). Finally multi-class SVM is non-trivial as SVM is motivated by the geometry of binary classification.

⁴Note that unlike regression using MVN models, we cannot calculate this integral in closed form

Appendices

A Proofs

A.1 We prove the following from Section 1.3.1

Prove that when using the shared covariance matrix MVN model, the decision boundary is piecewise-linear:

The decision boundary for a class k' is given by the following set:

$$\{\mathbf{x} | p(y = k | \mathbf{x}; \hat{\mathbf{w}}) = p(y = k' | \mathbf{x}; \hat{\mathbf{w}})\}, \quad \forall k \neq k' \quad (15)$$

We can rewrite this set as follows:

$$\left\{ \mathbf{x} \mid \frac{p(\mathbf{x} | y = k; \hat{\mathbf{w}}) p(y = k)}{p(\mathbf{x} | y = k'; \hat{\mathbf{w}}) p(y = k')} = 1 \right\} \quad (16)$$

Which taking logs can be rewritten as:

$$\{\mathbf{x} | \log[p(\mathbf{x} | y = k; \hat{\mathbf{w}}) p(y = k)] - \log[p(\mathbf{x} | y = k'; \hat{\mathbf{w}}) p(y = k')] = 0\} \quad (17)$$

$$= \{\mathbf{x} | \log[N_{\mathbf{x}}(\hat{\boldsymbol{\mu}}_k, \boldsymbol{\Sigma})] + \log[p(y = k)] - \log[N_{\mathbf{x}}(\hat{\boldsymbol{\mu}}_{k'}, \boldsymbol{\Sigma})] - \log[p(y = k')] = 0\} \quad (18)$$

$$= \{\mathbf{x} | (\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k)^T \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \log(p(y = k)) \quad (19)$$

$$- [(\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k'})^T \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_{k'}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k'} + \log(p(y = k'))] = 0\} \quad (20)$$

We notice that the condition above simply says that two linear equations of \mathbf{x} are equal to each other. The \mathbf{x} that satisfy this must also therefore be described by a linear relationship.

A.2 Obtain the stated optimization problem in section 3 by finding the Lagrangian dual:

First we will define what the lagrangian dual is:

Definition A.1 Lagrangian Dual:

A technique that transforms a constrained problem into a less constrained problem. For a constrained problem:

$$\min_{\theta} f(\theta), \text{ subject to } g_i(\theta) \leq 0, \forall i \quad (21)$$

We can construct the lagrangian:

$$\text{mathcal{L}}(\boldsymbol{\lambda}) := \min_{\theta} f(\theta) + \sum_i \lambda_i g_i(\theta) \quad (22)$$

where $\lambda_i \geq 0$ are called **lagrangian multipliers**

We now will write the lagrangian dual for the soft-margin classifier:

$$\mathcal{L}(\boldsymbol{\lambda}, \tilde{\boldsymbol{\lambda}}) := \min_{\mathbf{w}, \epsilon} \|\mathbf{w}_1\|^2 + \sum_i \epsilon_i - \sum_i \lambda_i y_i (\langle \mathbf{w}_1, \mathbf{x}_i \rangle + w_0) + \lambda_i \epsilon_i - \lambda_i - \sum_i \lambda'_i \epsilon_i \quad (23)$$

with optimality conditions,

$$\mathbf{w}' = \frac{\sum_i \lambda_i y_i \mathbf{x}_i}{2}, \lambda_i + \lambda'_i = 1, \sum_i \lambda_i y_i = 0 \quad (24)$$

We will now rewrite this, first note that,

$$\|\mathbf{w}_1\|^2 = \frac{\|\sum_i \lambda_i y_i \mathbf{x}_i\|^2}{4} \quad (25)$$

And that,

$$- \sum_i \lambda_i y_i \langle \mathbf{w}_1, \mathbf{x}_i \rangle = - \frac{\langle \sum_i \lambda_i y_i \mathbf{x}_i, \sum_i \lambda_i y_i \mathbf{x}_i \rangle}{2} \quad (26)$$

and that eq. (25) + eq. (26) is,

$$-\frac{\langle \sum_i \lambda_i y_i \mathbf{x}_i, \sum_i \lambda_i y_i \mathbf{x}_i \rangle}{4} = -\frac{\tilde{\boldsymbol{\lambda}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\lambda}}}{4} \quad (27)$$

where $\tilde{\boldsymbol{\lambda}} = [\lambda_1 y_1 \cdots \lambda_n y_n]^T$. Also note that

$$\sum_i \epsilon_i - \sum_i \lambda_i \epsilon_i - \sum_i \lambda'_i \epsilon_i \quad (28)$$

$$= \sum_i \epsilon_i (1 - \lambda_i - \lambda'_i) \quad (29)$$

$$= 0 \quad (30)$$

finally note,

$$-\sum_i \lambda_i y_i w_0 = -w_0 \sum_i \lambda_i y_i = 0 \quad (31)$$

Note that $\mathcal{L}(\boldsymbol{\lambda}, \tilde{\boldsymbol{\lambda}}) = \text{eq. (25)} + \text{eq. (26)} + \text{eq. (28)} + \text{eq. (31)} + \sum_i \lambda_i$, therefore,

$$\mathcal{L}(\boldsymbol{\lambda}, \tilde{\boldsymbol{\lambda}}) := -\frac{\tilde{\boldsymbol{\lambda}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\lambda}}}{4} + \sum_i \lambda_i \quad (32)$$

So our lagrangian problem is:

$$\max_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}) := \max_{\boldsymbol{\lambda}} -\frac{\tilde{\boldsymbol{\lambda}} \mathbf{X}^T \mathbf{X} \boldsymbol{\lambda}}{4} + \langle \tilde{\boldsymbol{\lambda}}, \mathbf{1} \rangle \quad (33)$$

Under constraints:

$$0 \leq \lambda_i \leq 1, \sum_i \lambda_i y_i = 0 \quad (34)$$

where $0 \leq \lambda_i \leq 1$ comes from the fact that,

$$\forall_i : \lambda_i \geq 0, \lambda'_i \geq 0, \text{ and } \lambda_i + \lambda'_i = 1 \quad (35)$$

B Homeworks

B.1 For Section 1

Question B.1.1 *Proof in section A.1*

Question B.1.2 *Derive the MLE for the parameters of the multinomial distribution:*
Let $\mathbf{p} = [p_1, \dots, p_d]^T$, where p_i denotes the probability of event i occurring and d is the number of mutually exclusive events. Then we can write out the likelihood:

$$p(D|\mathbf{p}) = \sum_{j=1}^n \frac{k!}{x_j^{(1)} \cdots x_j^{(d)}} \prod_{i=1}^d p_i^{x_j^{(i)}} = k! \sum_{j=1}^n \prod_{i=1}^d \frac{p_i^{x_j^{(i)}}}{x_j^{(i)}!} \quad (36)$$

Now taking the log we get:

$$\log(p(D|\mathbf{p})) = \log(k!) + \sum_{j=1}^n \sum_{i=1}^d x_j^{(i)} \log(p_i) - \log(x_j^{(i)}!) \quad (37)$$

We must have that $\sum_{i=1}^d p_i = 1$, the lagrangian with this constraint is:

$$\mathcal{L}(\mathbf{p}, \lambda) = \log(p(D|\mathbf{p})) + \lambda(1 - \sum_{i=1}^d p_i) \quad (38)$$

To find the MLE we now differentiate the lagrangian wrt. p_i :

$$\frac{\partial}{\partial p_i} \mathcal{L}(\mathbf{p}, \lambda) = \frac{\sum_{j=1}^n x_j^{(i)}}{\sum_{j=1}^n p_i} - \lambda \quad (39)$$

Setting this equal to zero and solving we get:

$$p_i = \sum_{j=1}^n \frac{x_j^{(i)}}{\lambda} \quad (40)$$

Now using our constraint we solve for λ :

$$\sum_{i=1}^d p_i = \sum_{i=1}^d \sum_{j=1}^n \frac{x_j^{(i)}}{\lambda} \quad (41)$$

$$\Rightarrow 1 = \frac{1}{\lambda} \sum_{i=1}^d \sum_{j=1}^n x_j^{(i)} \quad (42)$$

$$\Rightarrow \lambda = \sum_{j=1}^n \sum_{i=1}^d x_j^{(i)} \quad (43)$$

Therefore our MLE of the parameters is:

$$\hat{p}_i = \frac{\sum_{j=1}^n x_j^{(i)}}{\sum_{j=1}^n \sum_{i=1}^d x_j^{(i)}} \quad (44)$$

Question B.1.3 *Explain the Naive Bayes classifier using a maximum likelihood framework:*

Lets say we would like to model a classifier using the multinomial distribution as our model. When choosing to use the multinomial distribution as our model, we "naively" assume that all the features are mutually independent. We would like to formulate a prediction:

$$\hat{y} := p(y|\mathbf{x}; \mathbf{p}) \quad (45)$$

Using Bayes rule we have that:

$$p(y|D, \mathbf{p}) = \frac{p(D|\mathbf{p})p(y)}{p(D)} \quad (46)$$

$$\propto p(y) \cdot p(D|\mathbf{p}) \quad (47)$$

Using the previous question we can find using the MLE a $\hat{\mathbf{p}}$ that maximizes $p(D|\mathbf{p})$, then $p(y)$ can simply be calculated by finding the proportion of the time in the dataset where y equals a certain value. This then allows us to formulate a prediction, that has the maximum likelihood, for the classification of a given input \mathbf{x} as long as our "naive" assumptions are true.

B.2 For Section 2

Question B.2.1 What are the decision functions given by a binary logistic regression? One of them is $p(y|\mathbf{x}; \mathbf{w}) - \frac{1}{2}$ the other is $\frac{1}{2} - p(y|\mathbf{x}; \mathbf{w})$.

Question B.2.2 Prove that if $p(\mathbf{x}|y = 1)$ and $p(\mathbf{x}|y = -1)$ are MVN with shared covariance matrix Σ but different means μ_+, μ_- :

Sub-Question B.2.2.1 $\exists \mathbf{w}^*$ such that $p(y|\mathbf{x}) = \sigma((\langle \mathbf{x}, \mathbf{w}_1^* \rangle + w_0^*) \cdot y)$:

We have that,

$$p(y|\mathbf{x}; \mathbf{w}) = \sigma(f(\mathbf{x}; \mathbf{w}) \cdot y) \quad (48)$$

and that,

$$f(\mathbf{x}; \mathbf{w}) = \log \left(\frac{p(\mathbf{x}|y = 1)}{p(\mathbf{x}|y = -1)} \right) \quad (49)$$

Note that the log ratio of these densities makes up the binary decision boundary for this classification problem. Also recall that in section 1.3.1 we claimed that if $p(\mathbf{x}|y = k; \mathbf{w}) = N_{\mathbf{x}}(\mu_k, \Sigma)$ for all classes k (for some covariance matrix Σ), then the decision boundaries are piecewise linear, ie. in the binary case $\exists \mathbf{w}^*$ such that:

$$f(\mathbf{x}, \mathbf{w}) = \langle \mathbf{x}, \mathbf{w}_1^* \rangle + w_0^* \quad (50)$$

hence,

$$p(y|\mathbf{x}) = \sigma((\langle \mathbf{x}, \mathbf{w}_1^* \rangle + w_0^*) \cdot y) \quad (51)$$

Sub-Question B.2.2.2 Find \mathbf{w}^* :

$$f(\mathbf{x}; \mathbf{w}) = \log \left(\frac{p(\mathbf{x}|y = 1)}{p(\mathbf{x}|y = -1)} \right) \quad (52)$$

$$= \log \left(\frac{N_{\mathbf{x}}(\mu_+, \Sigma)}{N_{\mathbf{x}}(\mu_-, \Sigma)} \right) \quad (53)$$

$$= \log \left(\frac{(2\pi)^{-d/2} \det(\Sigma)^{-1/2} \exp(-\frac{1}{2}(\mathbf{x} - \mu_+)^T \Sigma^{-1}(\mathbf{x} - \mu_+))}{(2\pi)^{-d/2} \det(\Sigma)^{-1/2} \exp(-\frac{1}{2}(\mathbf{x} - \mu_-)^T \Sigma^{-1}(\mathbf{x} - \mu_-))} \right) \quad (54)$$

$$= -\frac{1}{2}(\mathbf{x} - \mu_+)^T \Sigma^{-1}(\mathbf{x} - \mu_+) + \frac{1}{2}(\mathbf{x} - \mu_-)^T \Sigma^{-1}(\mathbf{x} - \mu_-) \quad (55)$$

$$= (\Sigma^{-1} \mu_+)^T \mathbf{x} - \frac{1}{2} \mu_+^T \Sigma^{-1} \mu_+ - (\Sigma^{-1} \mu_-)^T \mathbf{x} + \frac{1}{2} \mu_-^T \Sigma^{-1} \mu_- \quad (56)$$

$$= ((\Sigma^{-1} \mu_+)^T - (\Sigma^{-1} \mu_-)^T) \mathbf{x} + \left(\frac{1}{2} \mu_-^T \Sigma^{-1} \mu_- - \frac{1}{2} \mu_+^T \Sigma^{-1} \mu_+ \right) \quad (57)$$

Hence we have found \mathbf{w}^* .

Question B.2.3 What is the probabilistic interpretation of f ?:

First, continuing from the end of section 2 we will describe using multi-class logistic regression with one-hot encoding. Note:

$$f(\mathbf{x}; \mathbf{w}) = \mathbf{W}^T \tilde{\mathbf{x}}, \mathbf{w} \in \mathbb{R}^{d \times K}, \tilde{\mathbf{x}} := \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \quad (58)$$

For $y_i \in \{1, \dots, K\}$ we let $\mathbf{t}_i \in \mathbb{R}^K$ be the one-hot encoding for it, further let:

$$\sigma(f, t) := \frac{\exp\langle f, t \rangle}{\sum_k \exp f^{(k)}} \quad (59)$$

Then:

$$\mathbf{w}_{MLE} = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i; \mathbf{w}), \mathbf{t}_i) \quad (60)$$

In this scenario what's the probabilistic interpretation of f ? like from which we derived the model of f in the binary logistic regression case in section 2. Recall that in discriminative classifiers we are trying to infer $p(y|\mathbf{x}; \mathbf{w})$ and that:

$$p(y|\mathbf{x}; \mathbf{w}) = \frac{p(\mathbf{x}|y; \mathbf{w})p(y)}{\sum_{y'} p(\mathbf{x}|y'; \mathbf{w})p(y')} \quad (61)$$

Let $i \in \{1, \dots, K\}$, consider $p(y = i|\mathbf{x}; \mathbf{w})$,

$$p(y = i|\mathbf{x}; \mathbf{w}) = \frac{p(\mathbf{x}|y = i; \mathbf{w})p(y = i)}{\sum_{y'} p(\mathbf{x}|y'; \mathbf{w})p(y')} \quad (62)$$

$$= \frac{1}{1 + \frac{\sum_{y', y' \neq i} p(\mathbf{x}|y'; \mathbf{w})p(y')}{p(\mathbf{x}|y=i; \mathbf{w})p(y=i)}} \quad (63)$$

Like in the binary case we can focus on modelling the density ratio, so the probabilistic interpretation of f is:

$$f(\mathbf{x}; \mathbf{w})^{(i)} = \log \left(\frac{p(\mathbf{x}|y = i; \mathbf{w})p(y = i)}{\sum_{y', y' \neq i} p(\mathbf{x}|y'; \mathbf{w})p(y')} \right) \quad (64)$$

We further note that if $\hat{y} := \operatorname{argmax}_y p(y|\mathbf{x}; \mathbf{w})$, then our predictions will correspond to the multi-class decision rule in definition 1.2, as:

$$\operatorname{argmax}_y p(y|\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_y f(\mathbf{x}; \mathbf{w})^{(y)} \quad (65)$$