

DS-GA 1003 Machine Learning: Homework 1

Due 11.59 p.m. EST, February 27, 2024 on Gradescope

(fill in your name here)
(collaborators if any)

We encourage **L^AT_EX**-typeset submissions but will accept quality scans of hand-written pages.

1 Linear Regression Model

Consider the data generating process as such: $\mathbf{x} \in \mathbb{R}^D$ is drawn from some unknown $p(\mathbf{x})$ and $y = w_1^{\text{true}} x_1 + \epsilon_y$, where $w_1^{\text{true}} \in \mathbb{R}$ and $\epsilon_y \sim \mathcal{N}(0, 1)$. This is unknown to us, as a result, we construct a linear model for y using all D features of \mathbf{x} , instead of just using x_1 .

- (A) Explain what the terms **model class** and **model misspecification** mean. Is our model correctly *specified* here? Why or why not?

Solution. The model class is linear regression model and is correctly chosen since it captures the real data generation process (which is also a linear regression model).

However the model with w_1^{all} contains some model misspecification. Also it is a linear model, we choose all D features of the data to fit the model, which is different from data generation process, which only uses the first dimension of the feature. This could cause multicollinearity problem as we will explore below. \square

Let \widehat{w}_1 be the estimate of w_1^{true} using only x_1 and let $\widehat{w}_1^{\text{all}}$ be the estimate of w_1^{true} when using all of \mathbf{x} . We will study the effects of our model by analyzing the relationships between $\mathbb{E}[\widehat{w}_1^{\text{all}}]$ and $\mathbb{E}[\widehat{w}_1]$, as well as between $\text{Var}[\widehat{w}_1^{\text{all}}]$ and $\text{Var}[\widehat{w}_1]$. We do so empirically by running PyTorch simulations as follows:

1. Pick any value of w_1^{true} you like as ground truth, e.g. with `torch.randn(1)`.
2. Write a function, taking D and c as input, that does the following: **(1)** Generate $N = 50$ samples of $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, where Σ is the $D \times D$ covariance matrix with all diagonal entries equal to $\sigma^2 = 1$ and all off-diagonal entries equal to c . **(2)** Compute y using the relationship above (note that y only depends on the first feature). This involves drawing N samples of noise $\epsilon_y \sim \mathcal{N}(0, 1)$. **(3)** Using our dataset of N samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, compute the least-squares solutions for $\widehat{w}_1^{\text{all}}$ (i.e. using all features and taking the first coefficient) and \widehat{w}_1 (i.e. using only one feature).
3. Write a function that performs Step 2 for $T = 100$ trials, i.e. each trial generates a new dataset to compute $\widehat{w}_1^{\text{all}}$ and \widehat{w}_1 . (Note that w_1^{true} is constant throughout.) For each estimator, compute the mean and standard deviation of the T trials.
4. Perform Step 3 for each $c \in \{0.1, 0.2, 0.3, \dots, 0.9\}$ and each $D \in \{2, 4, 8, 16, 32\}$. Separately plot the means and standard deviations as a function of c , using the same plot for both estimators. This means you should have 10 plots altogether: two plots (means and standard deviations) for each of the five choices of D . Each plot will contain two curves (the two estimators).

- (B) What do you observe with respect to c and D ? How do you explain your results? Show a few (not all) of your 10 generated plots to support your answer.

Solution. The observations are as follows:

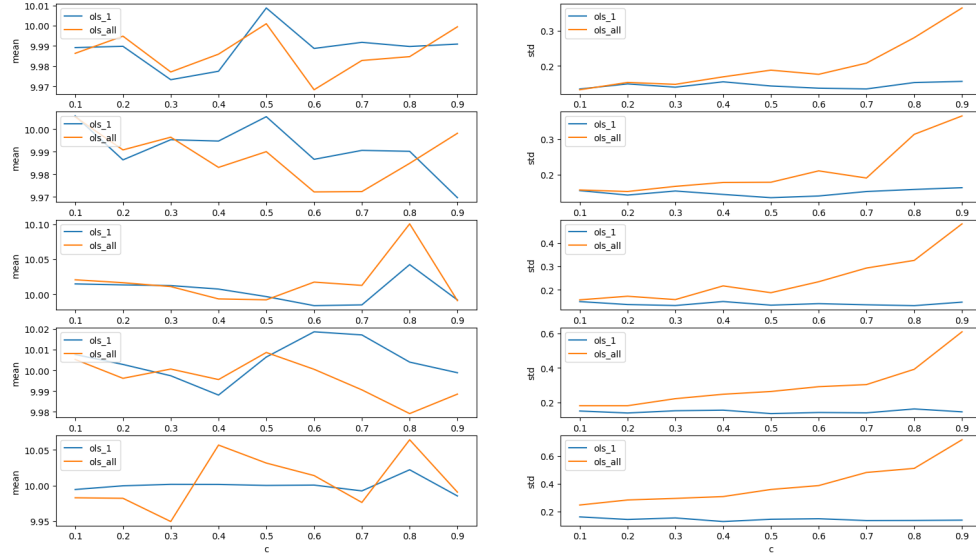


Figure 1: Grid Look

- (a) Observation 1: As c increases (feature correlation), we have higher variance in \hat{w}_1^{all} and stable variance in \hat{w}_1 .
- (b) Observation 2: As c increases (feature correlation), we have both \hat{w}_1^{all} and \hat{w}_1 around w_{true} .
- (c) Observation 3: As D increases, we don't observe any strange behaviors across c .

Then we give a detailed proof for each of these observations.

Claim 1: The \hat{w}_1 is both unbiased and consistent estimator for w_{true}

Proof. For \hat{w}_1 , our linear model is defined as:

$$y_i = \omega_1 x_i + \varepsilon_y$$

where $\varepsilon_y \sim N(0, 1)$ The OLS estimator given dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, is:

$$\hat{w}_1 = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2}$$

We first show that $\hat{w}_1 \xrightarrow{P} \omega_1^{true}$ Since $Y = \omega_1^{true} X + \varepsilon_y$, we have:

$$\begin{aligned} \text{cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &= E[W_1^{true} X^2 + X\varepsilon_y] - 0 \\ &= \omega_1^{true} E[X^2] + E[X\varepsilon_y] \\ &= \omega_1^{true} \\ \text{Var}[X] &= 1 \end{aligned}$$

$$\frac{1}{N-1} \sum_{i=1}^N x_i y_i \xrightarrow{P} \text{Cov}(X, Y)$$

$$\frac{1}{N-1} \sum_{i=1}^N x_i^2 \xrightarrow{P} \text{Var}(X)$$

by WLLN.

\therefore By continuous mapping theorem.

$$\frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2} \xrightarrow{P} \frac{\text{Cov}(X, Y)}{\text{Var}(x)} = \omega_1^{\text{true}}$$

We then show $E[\hat{w}_1] = \omega_1^{\text{true}}$.

$$\begin{aligned} \therefore E[\hat{w}_1] &= \frac{E\left[\sum_{i=1}^N x_i (\omega_1^{\text{true}} x_i + \varepsilon_y)\right]}{\sum_{i=1}^N x_i^2} \\ &= \omega_1^{\text{true}} + \frac{E\left[\sum_{i=1}^N x_i \varepsilon_y\right]}{\sum_{i=1}^N x_i^2} \\ &= \omega_1^{\text{true}} \end{aligned}$$

□

Claim 2: The $\widehat{w_1^{\text{all}}}$ is unbiased but not a consistent estimator for w_1^{true}

Proof. For $\widehat{w_1^{\text{all}}}$, we defined our linear model as:

$$\vec{y} = X\vec{w} + \vec{\varepsilon}$$

where. $\widehat{w_1^{\text{all}}} = \vec{w}_1^*$, $\vec{\varepsilon} \sim N(0, I_N)$

The OLS estimator is given by:

$$\vec{w}^* = \left(X^\top X\right)^{-1} X^\top \vec{y}$$

We show that $E[\vec{w}^*] = \vec{w}^{\text{true}}$

$$\begin{aligned} \therefore E[\vec{w}^*] &= E\left[\left(X^\top X\right)^{-1} X^\top (X\vec{w}^{\text{true}} + \vec{\varepsilon})\right] \\ &= \left(X^\top X\right)^{-1} \left(X^\top X\right) \vec{w}^{\text{true}} + E\left[\left(X^\top X\right)^{-1} X^\top \vec{\varepsilon}\right] \\ &= \vec{w}^{\text{true}} + \vec{0} \\ &= \vec{w}^{\text{true}} \\ \therefore E[\vec{w}_i^*] &= \omega_i^{\text{true}} \forall i = 1, 2, \dots, D \end{aligned}$$

We then show $\text{Var}[\vec{w}^*] = \left(X^\top X\right)^{-1}$ Since $\text{Var}(\vec{w}^*) = E\left[(\vec{w}^* - \vec{w}^{\text{true}})(\vec{w}^* - \vec{w}^{\text{true}})^\top\right]$

$$\begin{aligned} &= E\left[\left(X^\top X\right)^{-1} X^\top \vec{\varepsilon} \vec{\varepsilon}^\top X \left(X^\top X\right)^{-1}\right] \\ &= \left(X^\top X\right)^{-1} X^\top E\left[\vec{\varepsilon} \vec{\varepsilon}^\top\right] X \left(X^\top X\right)^{-1} \\ &= \left(X^\top X\right)^{-1} X^\top I_N X \left(X^\top X\right)^{-1} \\ &= \left(X^\top X\right)^{-1} \end{aligned}$$

□

\therefore When c is large, there's strong or perfect multicollinearity in data matrix X , which means if \vec{v} is an eigenvector of $X^\top X$ then $X^\top X \vec{v} = X^\top (\vec{0}) = \vec{0}$. since some columns can be linear combination of other columns.

\therefore We could have small eigenvalues, close to zero. Since $\lambda_i(X^\top X) = \lambda_i\left((X^\top X)^{-1}\right)^{-1}$

\therefore We expect large eigenvalues from $(X^\top X)^{-1}$, which corresponds to the variance of the entries of \hat{w}^{all} , thus we observe large variance of w_1^{all} at larger C values. \square

2 Bayesian Linear Regression Model

Consider the data generating process as such: $x \sim \mathcal{N}(0, 1)$ and $y = w_{true}x^2 + \epsilon$, where $w_{true} = 1.0$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$, and $\sigma^2 = 1.0$. Again, this is unknown to us. We will model the data using Bayesian linear regression.

- (A) Using PyTorch, simulate a dataset $\mathcal{D}_N = \{(x_i, y_i)\}_{i=1}^N$ for each $N \in \{10, 100, 1000, 10000\}$ according to the true data generating process above. For our Bayesian linear regression model, let us choose our prior as $w \sim N(0, 1)$ and our likelihood as $y|w, x \sim N(wx, \sigma^2)$ for $\sigma^2 = 1.0$. Compute the mean and variance of the posterior $w|\mathcal{D}_N$ for each dataset. Does the posterior concentrate on w_{true} ? Why or why not?

Solution. Suppose we have dataset $D_N = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, where we assume a 1-d model $y_i = \omega x_i + \varepsilon_i$ where $\varepsilon_i \sim N(0, \sigma^2)$, $\omega \sim N(\mu_0, \sigma_0^2)$

Then the posterior distribution:

$$\begin{aligned}
 f(\omega | D_N) &= f(\omega | x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N, \sigma^2) \\
 &= \frac{f(y_1, y_2, \dots, y_N | \omega, x_1, x_2, \dots, x_N, \sigma^2) f(\omega)}{f(y_1, y_2, \dots, y_N | x_1, x_2, \dots, x_N)} \\
 &\propto \prod_{i=1}^N \exp\left\{-\frac{1}{2\sigma^2} (y_i - \omega x_i)^2\right\} \exp\left\{-\frac{1}{2\sigma_0^2} (\omega - \mu_0)^2\right\} \\
 &\propto -\frac{\sigma_0^2 \sum_{i=1}^N (y_i - \omega x_i)^2 + \sigma^2 (\omega - \mu_0)^2}{2\sigma^2 \sigma_0^2} \\
 &= -\frac{\left(\sum_{i=1}^N x_i^2 \cdot \omega^2 - 2 \sum_{i=1}^N x_i y_i \cdot \omega + \sum_{i=1}^N y_i^2\right) \sigma_0^2 + (\omega^2 - 2\mu_0 \cdot \omega + \mu_0^2) \cdot \sigma^2}{2\sigma^2 \sigma_0^2} \\
 &= -\frac{\left(\sigma_0^2 \sum_{i=1}^N x_i^2 + \sigma^2\right) \omega^2 - 2 \left(\sigma_0^2 \sum_{i=1}^N x_i y_i + \mu_0 \cdot \sigma^2\right) \omega + C}{2\sigma^2 \sigma_0^2} \\
 \therefore \mu_{\text{post}} &= \frac{\sigma_0^2 \sum_{i=1}^N x_i y_i + \mu_0 \sigma^2}{\sigma_0^2 \sum_{i=1}^N x_i^2 + \sigma^2}, \sigma_{\text{post}}^2 = \frac{\sigma^2 \sigma_0^2}{\sigma_0^2 \sum_{i=1}^N x_i^2 + \sigma^2}
 \end{aligned}$$

, where in this question we have $\sigma_0^2 = \sigma_1^2 = 1, \mu_0 = 0$, and that as we notice from the graph $\mu_{\text{post}} = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2 + 1}$, $\sigma_{\text{post}}^2 = \frac{1}{\sum_{i=1}^N x_i^2 + 1}$. As shown in the graph:

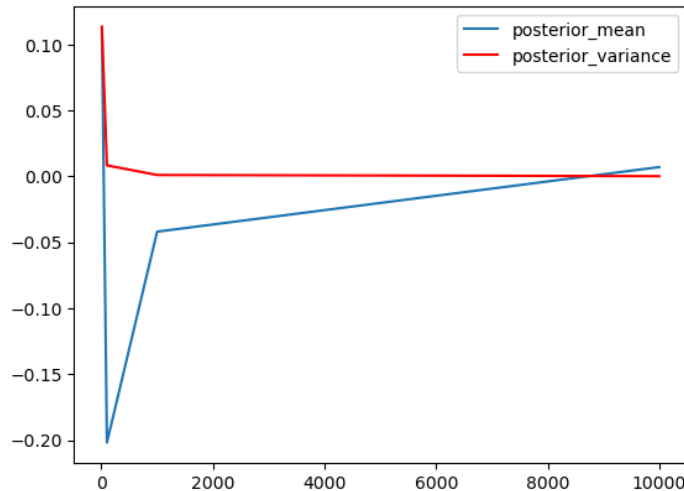


Figure 2: Main name 2

the posterior variance is decreasing as we increasing the number of data points, this is the same to increase the denominator of posterior variance. The posterior mean will not concentrate around w_{true} since $x \sim \mathcal{N}(0, 1)$, which makes the $\sum_{i=1}^n x_i y_i$ term centered around 0 when $n \rightarrow \infty$. \square

- (B) What would be challenging about our analysis in part (A) if we had picked a different prior, for example, Laplace or Gamma?

Solution. If we pick different prior, it may be difficult to compute the posterior. For example, for Gamma prior, it is not a conjugate prior for gaussian likelihood (but is for poisson likelihood), so the posterior is not Gamma. The same happens for Laplacian prior, the posterior is not gaussian either. But if we use MAP estimate, then even if we don't know the exact posterior distribution, we may still get our estimation \hat{w} . \square

- (C) Repeat part (A), except we use the basis set $\phi(x) = [x, x^2]$ (instead of x itself) and perform 2D Bayesian linear regression. We choose our prior to be $\mathbf{w} \sim N(\mathbf{0}, \mathbf{I})$, where \mathbf{I} is the 2×2 identity matrix, and our likelihood to be $y|\mathbf{w}, x \sim N(\mathbf{w}^\top \phi(x), \sigma^2)$ for $\sigma^2 = 1.0$. Compute the mean and variance of the posterior $\mathbf{w}|\mathcal{D}_N$ for each dataset. What do you observe about the posterior as N changes? Why?

Solution. The statistics are summarized in the following table:

N	mean	covariance matrix
10	[0.2046, 0.3908]	$\begin{bmatrix} 0.0749 & -0.0020 \\ -0.0020 & 0.0330 \end{bmatrix}$
100	[-0.0693, 0.2350]	$\begin{bmatrix} 0.0079 & -0.0004 \\ -0.0004 & 0.0021 \end{bmatrix}$
1000	[0.0149, 0.3488]	$\begin{bmatrix} 1.0616 \times 10^{-3} & 5.5125 \times 10^{-6} \\ 5.5125 \times 10^{-6} & 4.1348 \times 10^{-4} \end{bmatrix}$
10000	[0.0149, 0.3488]	$\begin{bmatrix} 9.9178 \times 10^{-5} & -3.6099 \times 10^{-7} \\ -3.6099 \times 10^{-7} & 3.2670 \times 10^{-5} \end{bmatrix}$

We notice that as $N \rightarrow \infty$, the variance of the posterior distribution converges to zero and the mean of the posterior distribution could be anywhere, but generally, it is the weighted average between the likelihood mean and the prior mean. \square

- (D) Reflecting on your answers in parts (A) and (C), name one challenge that **cannot be solved** by using a Bayesian model (instead of a frequentist approach like standard linear regression).

Solution. One challenge is the choice of prior distribution. If our prior variance is too small and the prior mean is too far away from the true parameter, then small prior variance causes the posterior distribution to put more weight on the prior during each bayesian updating iteration. Moreover, if we choose a bad prior that has no conjugacy with its likelihood, the posterior distribution may be hard to compute. \square

- (E) Reflecting on your answers in part (C) and in Question 1, name one challenge that **can be improved** by using a Bayesian model.

Solution. One improvement with bayesian model is that our posterior variance for our estimator will always be smaller for each iteration of bayesian updating. In the case of multicollinearity features, unregularized linear regression model like in problem 1 will generally overfit to the dataset and have larger variance in the estimator while for bayesian model, adding prior is the same to adding the regularization terms to the parameter, making the variance of the estimator within bounds.

For example, OLS(MLE) delivers $(X^\top X)^{-1}X^\top \vec{y}$ while MAP delivers $\vec{\mu}_W + (X^\top X + \Sigma_W^{-1})^{-1}X^\top (\vec{y} - X\vec{\mu}_W)$ where $\vec{\omega} \sim \mathcal{N}(\mu_W, \Sigma_W)$ (prior distribution). We could easily see that the MAP results guarantee the invertibility of $X^\top X$ while fall victims to the choice of prior distribution where $\vec{\mu}_W$ and Σ_W could greatly impact the location of MAP estimate. \square

Hint: Both part (C) above and Question 1 involve doing linear regression with many correlated features. How do these two sets of findings relate? Does using a Bayesian approach affect the way the model treats correlated features?

3 Class-Conditional Gaussian Generative Model

Consider a classification task where $\mathbf{x} \in \mathbb{R}^D$ and $y \in \{1, \dots, K\}$. We observe the dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$. Let us construct a model for the joint distribution as

$$p_\theta(\mathbf{x}, y) = p_\theta(\mathbf{x}|y)p_\theta(y)$$

where θ denotes the set of all parameters of the model.

- (A) Our model is known as a **class-conditional generative model**. What about the model makes it generative? What makes it class-conditional?

Solution. We break each term as follows:

- (a) Class-Conditional: This is because during the learning process, if we are giving the dataset \mathcal{D} , we can use MLE to fit a gaussian distribution with

$$\Sigma_k = \frac{1}{n_k} \sum_{i=1}^{n_k} (\vec{x}_i^{[k]} - \vec{\mu}_k)(\vec{x}_i^{[k]} - \vec{\mu}_k)^\top$$

$$\vec{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \vec{x}_i^{[k]}$$

for each class k such that we have a class-conditional probability distribution for the features $f(\vec{x}|y = k) \sim \mathcal{N}(\vec{\mu}_k, \Sigma_k)$

- (b) Generative: This is because given a class, we can generate(sample) new data points according to this class-conditional distribution.

□

- (B) For a given value of θ , how would you predict the label for a new test point \mathbf{x}_* using your model $p_\theta(\mathbf{x}, y)$?

Solution. Basically we will use the idea of MAP as follows:

$$y_* = \operatorname{argmax}_{y \in \{1, 2, \dots, k\}} P(y | \vec{x}_*)$$

$$= \operatorname{argmax}_{y \in \{1, 2, \dots, k\}} \frac{f(\vec{x}_* | y) P_\theta(y)}{\sum_{y \in \{1, \dots, k\}} f(\vec{x}_* | y) P_\theta(y)}$$

where $f(\vec{x}_* | y) \sim \mathcal{N}(\vec{\mu}_y, \Sigma_y)$.

□

Let us model y as a Categorical distribution $\text{Cat}(\boldsymbol{\pi})$. Here $p_\theta(y = k) \triangleq \pi_k$, where $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ such that $\forall k, \pi_k \geq 0$ and $\sum_k \pi_k = 1$. You may leave $p_\theta(\mathbf{x}|y)$ unspecified for now.

- (C) Write down an expression for the log-likelihood of the observed dataset \mathcal{D} .

Solution.

$$\mathcal{L}(\vec{\pi}) = \prod_{i=1}^N \prod_{j=1}^K \pi_j^{\mathbf{1}_{\{y_i=j\}}}$$

$$l(\vec{\pi}) = \sum_{i=1}^N \sum_{j=1}^K \mathbf{1}_{\{y_i=j\}} \log \pi_j$$

where $\sum_{j=1}^K \pi_j = 1$

□

- (D) Derive an expression for the maximum likelihood estimator (MLE) for $\boldsymbol{\pi}$, which we will denote as $\hat{\boldsymbol{\pi}}$. Make sure to account for the constraints on $\boldsymbol{\pi}$ using Lagrange multipliers.

Solution. The optimization problem becomes.

$$L(\boldsymbol{\pi}, \mu) = - \sum_{i=1}^N \sum_{j=1}^K \mathbf{1}\{y_i = j\} \log \pi_j + \mu \left(\sum_{j=1}^K \pi_j - 1 \right)$$

And we want:

$$\min_{\vec{\pi}} L(\vec{\pi}, \mu)$$

Now we prove that $\mathcal{L}(\vec{\pi}, \mu)$ is convex in $\vec{\pi}$.

$$\because \log(\pi_j) = \log(\vec{a}_j^\top \vec{\pi}) \text{ where } \vec{a}_j = [0, 0, \dots, 1, 0, 0]$$

$-\log(\vec{a}_j^\top \vec{\pi})$ is convex in $\vec{\pi}$ since $-\log x$ is convex in x and $-\log(\vec{a}_j^\top \vec{\pi})$ is a convex composition of a linear mapping. $\because -\sum_{i=1}^N \sum_{j=1}^K \mathbf{1}\{y_i = j\} \log \pi_j$ is nonnegative sum of convex functions and thus is convex.

$\because \mu \left(\sum_{j=1}^K \pi_j - 1 \right)$ is convex in $\vec{\pi}$ for similarly reason. $\therefore L(\vec{\pi}, \mu)$ is convex in $\vec{\pi}$

We enforce KKT conditions on $L(\vec{\pi}, \mu)$ guarantees optimality:

$$\begin{aligned} \frac{\partial \mathcal{L}(\vec{\pi}, \mu)}{\partial \pi_i} &= - \sum_{i=1}^N \frac{\mathbf{1}\{y_i = j\}}{\pi_j} + \mu = 0 \quad \forall j = 1, 2, \dots, K. \\ \sum_{i=1}^K \pi_i &= 1 \end{aligned}$$

Solving for the KKT we get:

$$\therefore \hat{\pi}_j = \frac{\sum_{i=1}^N \mathbf{1}\{y_i = j\}}{\mu} \quad \forall j = 1, 2, \dots, K.$$

We now solve for μ :

$$\begin{aligned} \because \sum_{j=1}^K \hat{\pi}_j &= 1 \\ \therefore \sum_{j=1}^K \sum_{i=1}^N \frac{\mathbf{1}\{y_i = j\}}{\mu} &= 1 \\ \therefore \mu &= \sum_{j=1}^K \sum_{i=1}^N \mathbf{1}\{y_i = j\} \text{ which is a constant} \end{aligned}$$

□

Let us further model $\mathbf{x}|y$ as (multivariate) Gaussian distributions $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \Sigma_k)$ for all K classes, where $\boldsymbol{\mu}_k \in \mathbb{R}^D$ and Σ_k is a $D \times D$ (positive semi-definite) covariance matrix. Assume that there are only $K = 2$ classes. This means that the total set of parameters are $\theta = \{\boldsymbol{\pi}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2\}$.

Now, consider the case where the data comes from this model. That is, $y \sim \text{Cat}(\boldsymbol{\pi}^{true})$ and $\mathbf{x}|y = k \sim \mathcal{N}(\boldsymbol{\mu}_k^{true}, \Sigma_k^{true})$ for all k . After we observe this data, we can then construct a *discriminative model* to predict y from \mathbf{x} , that is, we will learn a model for $p_{true}(y = k|\mathbf{x})$. Let us do so using **logistic regression**:

$$y|\mathbf{x} \sim \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x}))$$

where \mathbf{w} are the parameters of the model and $\sigma(z)$ is the logistic sigmoid:

$$\sigma(z) = \frac{1}{1 + \exp[-z]}$$

- (E) Will logistic regression always be able to model the true data conditional $p_{\text{true}}(y = k|\mathbf{x})$? If so, why? If sometimes, when? And if there are any cases where logistic regression will not be able to model $p_{\text{true}}(y = k|\mathbf{x})$, are there any ways to fix it?

Solution. The answer is not always. Let's say if logistic regression perfectly model the true data conditional $p_{\text{true}}(y = k | \mathbf{x})$, then since $K = 2$, we will have $\frac{p_{\text{true}}(y=1|\mathbf{x})}{p_{\text{true}}(y=0|\mathbf{x})} = \frac{1}{e^{-\mathbf{w}^\top \mathbf{x}}}$ and the decision boundary requires $\frac{1}{e^{-\mathbf{w}^\top \mathbf{x}}} = 1$ where $e^{\mathbf{w}^\top \mathbf{x}} = 1$ and $\mathbf{w}^\top \mathbf{x} = 0$, which is a linear boundary. But when $\Sigma_1 \neq \Sigma_2$, namely the covariance matrix from different class is different, the true decision boundary is quadratic instead of linear. The proof is as follows, we know in GDA, given test data point \mathbf{x}_* , our prediction for $\hat{\mathbf{y}}_*$ is calculated by MAP rules:

$$\begin{aligned} \text{MAP}(\mathbf{x}_*) &:= \arg \max_{y \in \{1, 2, \dots, K\}} p_{\tilde{y}|\tilde{\mathbf{x}}}(y | \mathbf{x}_*) \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} \frac{p_{\tilde{y}}(y) f_{\tilde{\mathbf{x}}|\tilde{y}}(\mathbf{x}_* | y)}{\sum_{k \in \{1, 2, \dots, k\}} p_{\tilde{y}}(k) p_{\tilde{\mathbf{x}}|\tilde{y}}(\mathbf{x}_* | k)} \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} \frac{\frac{\alpha_y}{\sqrt{(2\pi)^d |\Sigma_y|}} \exp\left(-\frac{1}{2} (\mathbf{x}_* - \mu_y)^T \Sigma_y^{-1} (\mathbf{x}_* - \mu_y)\right)}{\sum_{k \in \{1, 2, \dots, K\}} \frac{\alpha_k}{\sqrt{(2\pi)^d |\Sigma_k|}} \exp\left(-\frac{1}{2} (\mathbf{x}_* - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_* - \mu_k)\right)} \end{aligned} \quad (1)$$

and the decision boundary is determined by:

$$\frac{p_{\tilde{y}|\tilde{\mathbf{x}}}(y = 1 | \mathbf{x}_*)}{p_{\tilde{y}|\tilde{\mathbf{x}}}(y = 2 | \mathbf{x}_*)} = 1$$

using (1) we get:

$$\frac{\alpha_1 \sqrt{|\Sigma_2|}}{\alpha_2 \sqrt{|\Sigma_1|}} \exp\left(\frac{1}{2} (\mathbf{x}_* - \mu_2)^T \Sigma_2^{-1} (\mathbf{x}_* - \mu_2) - \frac{1}{2} (\mathbf{x}_* - \mu_1)^T \Sigma_1^{-1} (\mathbf{x}_* - \mu_1)\right) = 1$$

Then we take logarithms to get the decision boundary,

$$\frac{1}{2} (\mathbf{x}_* - \mu_2)^T \Sigma_2^{-1} (\mathbf{x}_* - \mu_2) - \frac{1}{2} (\mathbf{x}_* - \mu_1)^T \Sigma_1^{-1} (\mathbf{x}_* - \mu_1) + \log\left(\frac{\alpha_1 \sqrt{|\Sigma_2|}}{\alpha_2 \sqrt{|\Sigma_1|}}\right) = 0.$$

where if $\Sigma_1 = \Sigma_2$, we can cancel out the quadratic terms and obtain a linear boundary. Otherwise, the decision boundary is quadratic and the logistic regression fails to approximate the true distribution well. \square

4 Poisson Generalized Linear Model

Consider a classification task where $\mathbf{x} \in \mathbb{R}^D$ and $y \in \mathbb{N} = \{0, 1, 2, 3, \dots\}$, noting that the support of y is the unbounded set of natural numbers. We have an observed dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$. Let us also assume that the number of features, D , is larger than the number of examples, N . We will model this data using a Poisson Generalized Linear Model (GLM). Let $\boldsymbol{\theta}$ denote the linear coefficients of the model.

(A) Write down the log-likelihood function of the Poisson GLM.

Solution. Suppose we have dataset $D_N = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_N, y_N)\}$ where $\vec{x}_i \in \mathbb{R}^D$ the poisson model assumes, $f(y_i; \vec{x}_i) \stackrel{\text{ind}}{\sim} \text{poisson}(\lambda_i)$

$$\begin{aligned} \therefore f(y_i; \vec{x}_i) &= \frac{1}{y_i!} e^{-\lambda} \lambda^{y_i} \\ &= e^{y_i \cdot \log \lambda - \lambda} \cdot \frac{1}{y_i!} \\ &= e^{y_i \cdot \vec{x}_i^\top \vec{\theta} - e^{\vec{x}_i^\top \vec{\theta}}} f_0(y_i) \end{aligned}$$

where the link function is $g(\lambda) = \log \lambda$
The log-likelihood function is:

$$\begin{aligned} l(\vec{\theta}) &= \log \prod_{i=1}^N f(y_i; \vec{x}_i) \\ &= \log \exp \left\{ \sum_{i=1}^N y_i \cdot \left(\vec{x}_i^\top \vec{\theta} \right) - e^{\vec{x}_i^\top \vec{\theta}} \right\} + \log f_0(y_i) \\ &= \sum_{i=1}^N y_i \left(\vec{x}_i^\top \vec{\theta} \right) - e^{\vec{x}_i^\top \vec{\theta}} \end{aligned}$$

□

(B) Given a test point \mathbf{x}_\star and some estimate $\hat{\boldsymbol{\theta}}$ of the parameter, how do you make a prediction \hat{y}_\star ?

Solution. Given the link function $g(\lambda) = \log \lambda$, we know the prediction model:

$$\begin{aligned} \log \hat{y}_\star &= \mathbf{x}_\star^\top \hat{\boldsymbol{\theta}} \\ \hat{y}_\star &= e^{\mathbf{x}_\star^\top \hat{\boldsymbol{\theta}}} \end{aligned}$$

□

(C) Now suppose that the parameter $\hat{\boldsymbol{\theta}}$ of the Poisson GLM is estimated using ℓ_2 -regularized maximum likelihood estimation. If the test point \mathbf{x}_\star is *orthogonal* to the subspace generated by the training data, what is the distribution $\hat{y}_\star | \mathbf{x}_\star$ predicted by the Poisson GLM model? Prove your answer.

Solution.

□

(D) From your answer to part (C), motivate ℓ_1 -regularization when the number of features, D , is larger than the number of examples, N .

Solution.

□

5 Distances and Optimization Directions

Consider two pairs of distributions with mean and variance parameterization:

Pair 1: Normal(0, 0.0001), Normal(0.1, 0.0001)

Pair 2: Normal(0, 1000), Normal(0.1, 1000)

(A) Make two plots where each plot shows the pdfs for the distributions in the pair.

Solution. The distributions are shown below:

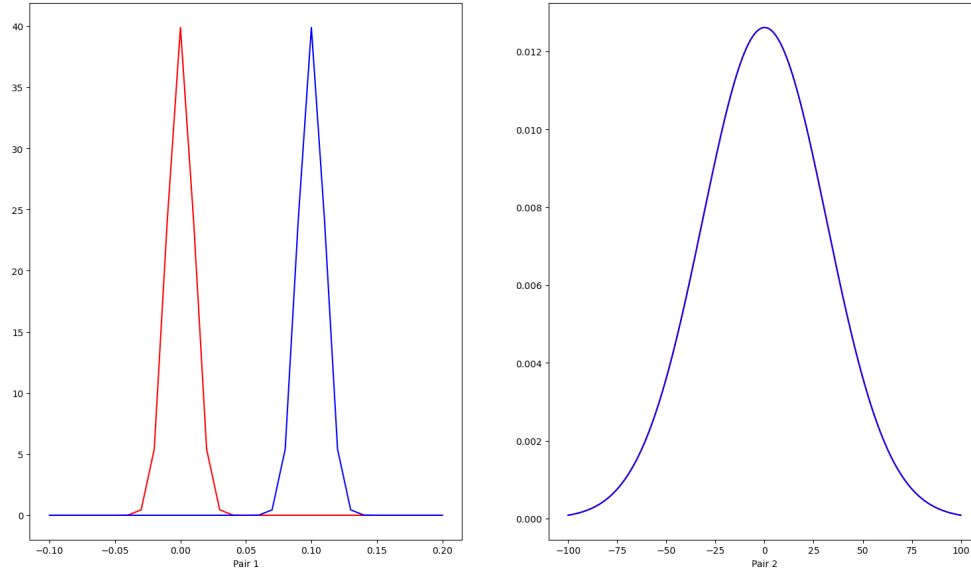


Figure 3: Two pairs of distribution

□

(B) Compute the Euclidean distance between the parameter vector (mean, variance) for both pairs of distributions. For the same pairs of distributions compute the KL-divergence. Which distance fits intuition better and why?

Solution. (a) The euclidean distance is given by:

$$\sqrt{(\mu_1 - \mu_2)^2 + (\sigma_1^2 - \sigma_2^2)^2}$$

using python we get:

i. $d(\text{pair1}) = 0.1$

ii. $d(\text{pair2}) = 0.1$

(b) The KL-Divergence is given by

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

using python we get:

i. $d(\text{pair1}) = 50.0$

ii. $d(\text{pair2}) = 5 \times 10^{-6}$

The KL divergence fits intuition better since it considers the actual probability distribution outcomes (spread) instead of just numerical values of mean and variance. In other words, for pair 2, due to larger variance, the outcomes from two distributions may overlap a lot, much like a low power test where the null distribution and alternative distribution overlaps a lot. \square

- (C) Assume θ_t is a parameter for a probability distribution and ρ_t is a scalar. What is the solution to the following optimization algorithm?

$$\max_{\theta_{t+1}} \sum_{i=1}^n \log p_{\theta_t}(y_i | \mathbf{x}_i) + (\theta_{t+1} - \theta_t)^\top \left[\nabla_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | \mathbf{x}_i) \Big|_{\theta=\theta_t} \right] - \frac{1}{2\rho_t} \|\theta_{t+1} - \theta_t\|_2^2$$

Solution. We immediately notice that it is a convex optimization problem w.r.t $\vec{\theta}_{t+1}$. Denote:

$$\mathcal{L}(\vec{\theta}_{t+1}) = \sum_{i=1}^n \log p_{\theta_t}(y_i | \mathbf{x}_i) + (\theta_{t+1} - \theta_t)^\top \left[\nabla_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | \mathbf{x}_i) \Big|_{\theta=\theta_t} \right] - \frac{1}{2\rho_t} \|\theta_{t+1} - \theta_t\|_2^2$$

, we have:

$$\begin{aligned} \nabla_{\vec{\theta}_{t+1}} \mathcal{L}(\vec{\theta}_{t+1}) &= -\frac{1}{\rho_t} (\vec{\theta}_{t+1} - \vec{\theta}_t) + \nabla_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | \mathbf{x}_i) \Big|_{\theta=\theta_t} = 0 \\ \nabla_{\vec{\theta}_{t+1}}^2 \mathcal{L}(\vec{\theta}_{t+1}) &= 2 \geq 0 \end{aligned} \tag{1}$$

thus the function is concave and we solve for (1) and could get: $\vec{\theta}_{t+1} = \vec{\theta}_t + \rho_t \nabla_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | \mathbf{x}_i) \Big|_{\theta=\theta_t}$ \square

- (D) What algorithm does the previous solution correspond to? Does part (B) say anything about why this algorithm might be suboptimal? How would you fix it?

Solution. The solution is just the gradient descent where the function we are descending is $\sum_{i=1}^n \log p_{\theta_t}(y_i | \mathbf{x}_i)$. Suboptimality results from the choice of distance metric between the parameters of two consecutive probability distributions $\vec{\theta}_{t+1}$ and $\vec{\theta}_t$. Here we are choosing euclidean distance as our distance metric where it will over-penalize the distance between two parameters. For example in problem (B) for if $\vec{\theta}_t$ is Normal(0, 1000) and $\vec{\theta}_{t+1}$ is Normal(0.1, 1000), then euclidean distance gives 0.1, which is huge compared with KL-divergence metric, which gives 5×10^{-6} . So a quick fix is to rewrite the optimization problem as follows:

$$\max_{\theta_{t+1}} \sum_{i=1}^n \log p_{\theta_t}(y_i | \mathbf{x}_i) + (\theta_{t+1} - \theta_t)^\top \left[\nabla_{\theta} \sum_{i=1}^n \log p_{\theta}(y_i | \mathbf{x}_i) \Big|_{\theta=\theta_t} \right] - \frac{1}{2\rho_t} D_{KL}(\theta_{t+1}, \theta_t)$$

which gives more accurate measurement of the distance. \square