

STA 314: Statistical Methods for Machine Learning I

Lecture 8 - Gradient Descent, Multi-class Logistic Regression

Xin Bing

Department of Statistical Sciences
University of Toronto

A general problem of solving a minimization problem

Suppose we want to solve the following problem

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \Theta} \mathcal{J}(\mathbf{w}; \mathcal{D}^{train}) := \operatorname{argmin}_{\mathbf{w} \in \Theta} \mathcal{J}(\mathbf{w})$$

where $\mathcal{J}(\mathbf{w}; \mathcal{D}^{train})$ is a differentiable function in \mathbf{w} , and depends on \mathcal{D}^{train} as well, and Θ is a subspace of \mathbb{R}^p .

- The optimal solution (if exists) must be a **critical point**, i.e. point to which the derivative is zero (partial derivatives to zero for multi-dimensional parameter).

Finding the optimal solution requires to solve the equations

- **Partial derivatives:** derivatives of a multivariate function with respect to one of its arguments.

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2) - f(x_1, x_2)}{h}$$

- The minimum must occur at a point where the partial derivatives are zero.

$$\begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{bmatrix} = 0$$

- This turns out to give a system of linear equations, which we can solve analytically in some scenarios.
- We may also use optimization techniques that iteratively get us closer to the solution.

- OLS:

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \mathcal{J}(\mathbf{w}) = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2.$$

The partial derivatives w.r.t. \mathbf{w} are

$$\frac{\partial g}{\partial \mathbf{w}} = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}).$$

(If not familiar with multi-dimensional derivatives, calculate $\frac{\partial g}{\partial w_j}$ and stack them together).

Setting the above equal to zero results

$$\mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}, \quad \Rightarrow \quad \hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- Ridge:

$$\hat{\mathbf{w}}_{\lambda}^R = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \mathcal{J}(\mathbf{w}) = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2.$$

The partial derivatives w.r.t. \mathbf{w} are

$$\frac{\partial g}{\partial \mathbf{w}} = -2\mathbf{X}^{\top}(\mathbf{y} - \mathbf{X}\mathbf{w}) + 2\lambda\mathbf{w}.$$

Setting the above equal to zero results

$$(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)\hat{\mathbf{w}}_{\lambda}^R = \mathbf{X}^{\top}\mathbf{y}, \quad \Rightarrow \quad \hat{\mathbf{w}}_{\lambda}^R = (\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^{\top}\mathbf{y}.$$

Gradient Descent

- Now let's see a second way to solve

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{J}(\mathbf{w})$$

which is more broadly applicable: **gradient descent**.

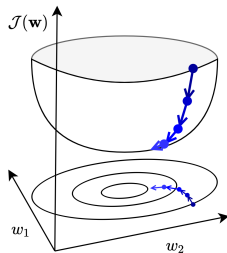
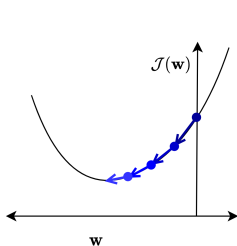
- Many times, we do not have a direct solution to

$$\frac{\partial \mathcal{J}}{\partial \mathbf{w}} = 0.$$

- Gradient descent is an **iterative algorithm**, which means we apply an update repeatedly until some criterion is met.

Gradient Descent

We **initialize** \mathbf{w} to something reasonable (e.g. all zeros) and repeatedly adjust them in the **direction of steepest descent**.



What is the direction of the steepest descent of $\mathcal{J}(\mathbf{w})$ at \mathbf{w} ?

Gradient Descent

- By definition, the direction of the greatest increase in $\mathcal{J}(\mathbf{w})$ at \mathbf{w} is its gradient $\partial\mathcal{J}/\partial\mathbf{w}$. So, we should update \mathbf{w} in the **opposite** direction of the gradient descent.
- The following update always decreases the cost function for small enough α (unless $\partial\mathcal{J}/\partial w_j = 0$): at the $(k + 1)$ th iteration,

$$w_j^{(k+1)} \leftarrow w_j^{(k)} - \alpha \cdot \frac{\partial\mathcal{J}}{\partial w_j} \Big|_{\mathbf{w}=\mathbf{w}^{(k)}}$$

- $\alpha > 0$ is a **learning rate** (or step size). The larger it is, the faster $\mathbf{w}^{(k+1)}$ changes relative to $\mathbf{w}^{(k)}$
 - ▶ We'll see later how to tune the learning rate, but values are typically small, e.g. 0.01 or 0.0001.

Example

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} \mathcal{J}(\mathbf{w}), \quad \mathcal{J}(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2.$$

Update rule in vector form at the $k + 1$ th iteration:

$$\begin{aligned} \mathbf{w}^{(k+1)} &\leftarrow \mathbf{w}^{(k)} - \alpha \frac{\partial \mathcal{J}}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^{(k)}} \\ &= \mathbf{w}^{(k)} + 2\alpha \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}^{(k)}). \end{aligned}$$

Initialization: $\mathbf{w}^{(0)} = \mathbf{0}$.

Stopping criteria

When do we stop?

- The objective value stops changing:

$$|\mathcal{J}(\mathbf{w}^{(k+1)}) - \mathcal{J}(\mathbf{w}^{(k)})| \text{ is small, i.e. } \leq 10^{-6}.$$

- The parameter stops changing: $\|\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}\|_2$ is small or $\|\mathbf{w}^{(k+1)} - \mathbf{w}^{(k)}\|_2 / \|\mathbf{w}^{(k)}\|_2$ is small.
- When we reach the maximum number (M) of iterations, e.g. $M = 1000$.

Gradient descent for solving the MLE under logistic regression

Recall we would like to solve

$$\min_{\mathbf{w} \in \mathbb{R}^p} \mathcal{J}(\mathbf{w})$$

where

$$\mathcal{J}(\mathbf{w}) = -\ell(\mathbf{w}) = \sum_{i=1}^n \left[-y_i \mathbf{x}_i^\top \mathbf{w} + \log \left(1 + e^{\mathbf{x}_i^\top \mathbf{w}} \right) \right].$$

The gradient at any \mathbf{w} is that, for any $j \in \{1, \dots, p\}$,

$$-\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_{i=1}^n \left[-y_i + \frac{e^{\mathbf{x}_i^\top \mathbf{w}}}{1 + e^{\mathbf{x}_i^\top \mathbf{w}}} \right] x_{ij} \quad (\text{verify this!})$$

Updates and stopping criteria

Therefore, at the $(k + 1)$ th iteration, with the learning rate α ,

$$\hat{\mathbf{w}}^{(k+1)} = \hat{\mathbf{w}}^{(k)} - \alpha \sum_{i=1}^n \left[-y_i + \frac{e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}}{1 + e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}} \right] \mathbf{x}_i.$$

Initialization $\mathbf{w}^{(0)} = \mathbf{0}$.

- The objective value stops changing: $|\ell(\hat{\mathbf{w}}^{(k+1)}) - \ell(\hat{\mathbf{w}}^{(k)})|$ is small, say, $\leq 10^{-6}$.
- The parameter stops changing: $\|\hat{\mathbf{w}}^{(k+1)} - \hat{\mathbf{w}}^{(k)}\|_2$ is small or $\|\hat{\mathbf{w}}^{(k+1)} - \hat{\mathbf{w}}^{(k)}\|_2 / \|\hat{\mathbf{w}}^{(k)}\|_2$ is small.
- Stop after M iterations for some specified M , e.g. $M = 1000$.

When should we expect Gradient Descent to work?

Recall we try to solve

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \Theta}{\operatorname{argmin}} \mathcal{J}(\mathbf{w}).$$

- Obviously, \mathcal{J} needs to be differentiable.
- If \mathcal{J} is also a convex function and Θ is a convex set, then Gradient Descent finds the optimal solution.
- In many cases, $\Theta = \mathbb{R}^p$ which is convex.

A set \mathcal{S} is convex if for any $\mathbf{x}_0, \mathbf{x}_1 \in \mathcal{S}$,

$$(1 - \lambda)\mathbf{x}_0 + \lambda\mathbf{x}_1 \in \mathcal{S} \quad \text{for all } 0 \leq \lambda \leq 1.$$

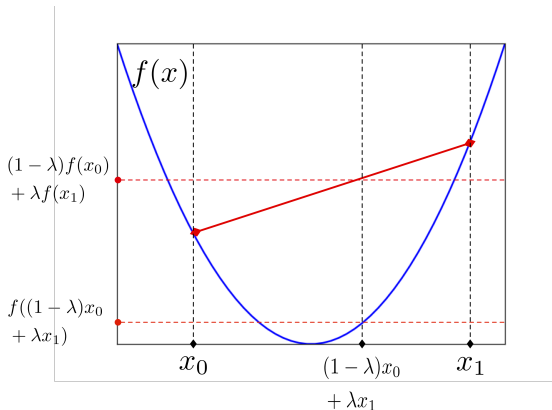
The Euclidean space \mathbb{R}^p is a convex set.

Convex Sets and Functions

- A function f is **convex** if for any $\mathbf{x}_0, \mathbf{x}_1$ in the domain of f ,

$$f((1 - \lambda)\mathbf{x}_0 + \lambda\mathbf{x}_1) \leq (1 - \lambda)f(\mathbf{x}_0) + \lambda f(\mathbf{x}_1), \quad \forall \lambda \in [0, 1].$$

- Equivalently, the set of points lying above the graph of f is convex.
- Intuitively: the function is bowl-shaped.



How to tell a loss is convex?

1. Verify the definition.
2. If f is twice differentiable and $f''(x) \geq 0$ for all x , then f is convex.
 - ▶ the least-squares loss function $(y - t)^2$ is convex as a function of t
 - ▶ the function
$$-yt + \log(1 + e^t)$$
is convex in t .
3. There are other sufficient conditions for convex, but non-differentiable, functions!

4 A composition rule: linear functions preserve convexity.

- ▶ If f is a convex function and g is a linear function, then both $f \circ g$ and $g \circ f$ are convex.
 - ▶ the least-square loss $(y - \mathbf{x}^\top \mathbf{w})^2$ is convex in \mathbf{w}
 - ▶ the negative log-likelihood under logistic regression

$$-y\mathbf{x}^\top \mathbf{w} + \log\left(1 + e^{\mathbf{x}^\top \mathbf{w}}\right)$$

is convex in \mathbf{w} .

- ▶ Both $\sum_i (y_i - \mathbf{x}_i^\top \mathbf{w})^2$ and $\sum_i \left[-y_i \mathbf{x}_i^\top \mathbf{w} + \log\left(1 + e^{\mathbf{x}_i^\top \mathbf{w}}\right) \right]$ are convex in \mathbf{w} .

There are more composition rules!

A great book:

Convex Optimization, Stephen Boyd and Lieven Vandenberghe.

Gradient Descent for Linear Regression

- The squared error loss

$$\sum_{i=1} (y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

of linear regression is a convex function. So there is a unique solution.

- Even in this case, we sometimes need to use GD.
- Why gradient descent, if we can find the optimum directly?
 - ▶ When p is large, GD is more efficient than direct solution
 - ▶ Linear regression solution: $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$
 - ▶ Matrix inversion is an $\mathcal{O}(p^3)$ algorithm
 - ▶ Each GD update costs $\mathcal{O}(np)$
 - ▶ Or less with stochastic GD (Stochastic GD, later)
 - ▶ Huge difference if $p \gg \sqrt{n}$

Gradient descent for solving the MLE under logistic regression

- The negative log-likelihood

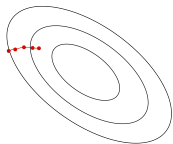
$$-\ell(\mathbf{w}) = \sum_{i=1}^n \left[-y_i \mathbf{x}_i^\top \mathbf{w} + \log \left(1 + e^{\mathbf{x}_i^\top \mathbf{w}} \right) \right]$$

is convex in \mathbf{w} .

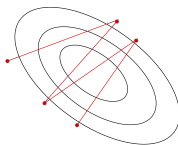
- So we can use gradient descent to find the minima of the logistic loss!
- GD can be applied to more general settings!

Effect of the Learning Rate (Step Size)

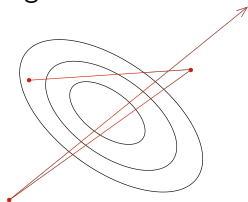
- In gradient descent, the learning rate α is a hyperparameter we need to tune. Here are some things that can go wrong:



α too small:
slow progress



α too large:
oscillations

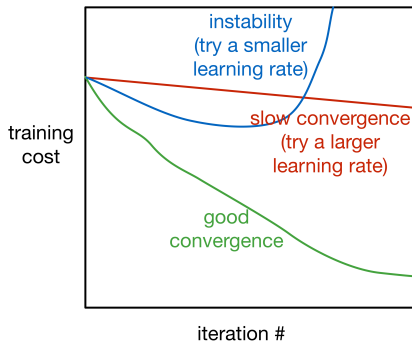


α much too large:
instability

- Good values are typically between 0.001 and 0.1. You should do a grid search if you want good performance (i.e. try 0.1, 0.03, 0.01, ...).

Training Curves

- To diagnose optimization problems, it's useful to look at the **training cost**: plot the training cost as a function of iteration.



- **Warning:** the training cost could be used to check whether the optimization problem reaches certain convergence. But
 - ▶ It does not tell whether we reach the global minimum or not
 - ▶ It does not tell anything on the performance of the fitted model

Gradient descent

Visualization:

http://www.cs.toronto.edu/~guerzhoy/321/lec/W01/linear_regression.pdf#page=21

Batch Gradient Descent

- Recall that

- ▶ OLS:

$$\hat{\mathbf{w}}^{(k+1)} = \hat{\mathbf{w}}^{(k)} + \alpha \sum_{i=1}^n \left[y_i - \mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)} \right] \mathbf{x}_i.$$

- ▶ Logistic regression:

$$\hat{\mathbf{w}}^{(k+1)} = \hat{\mathbf{w}}^{(k)} + \alpha \sum_{i=1}^n \left[y_i - \frac{e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}}{1 + e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}} \right] \mathbf{x}_i.$$

- Computing the gradient requires summing over **all** of the training examples, which can be done via matrix / vector operations. The fact that it uses all training samples is known as **batch training**.

Stochastic Gradient Descent

- Batch training is impractical if you have a large dataset (e.g. millions of training examples, $n \approx 10$ millions)!
- **Stochastic gradient descent (SGD)**: update the parameters based on the gradient for a single training example.

For each iteration $k \in \{1, 2, \dots\}$,

1. Choose $i \in \{1, \dots, n\}$ uniformly at random
2. Update the parameters by ONLY using this i th sample,

$$\begin{aligned}\hat{\mathbf{w}}^{(k+1)} &= \hat{\mathbf{w}}^{(k)} + \alpha \left[y_i - \mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)} \right] \mathbf{x}_i \\ \hat{\mathbf{w}}^{(k+1)} &= \hat{\mathbf{w}}^{(k)} + \alpha \left[y_i - \frac{e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}}{1 + e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}} \right] \mathbf{x}_i.\end{aligned}$$

Stochastic Gradient Descent

$$\hat{\mathbf{w}}^{(k+1)} = \hat{\mathbf{w}}^{(k)} + \alpha \left[y_i - \mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)} \right] \mathbf{x}_i$$
$$\hat{\mathbf{w}}^{(k+1)} = \hat{\mathbf{w}}^{(k)} + \alpha \left[y_i - \frac{e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}}{1 + e^{\mathbf{x}_i^\top \hat{\mathbf{w}}^{(k)}}} \right] \mathbf{x}_i.$$

Pros:

- Computational cost of each SGD update is independent of n !
- SGD can make significant progress before even seeing all the data!
- Mathematical justification: the gradients between SGD and GD have the same expectation for i.i.d. data.

Stochastic Gradient Descent

Cons: using single training example to estimate gradient:

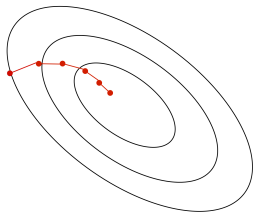
- Variance in the estimate may be high

Compromise approach:

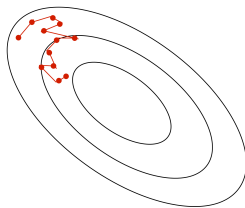
- compute the gradients on a randomly chosen medium-sized set of training examples $\mathcal{M} \subset \{1, \dots, n\}$, called a **mini-batch**.
- Stochastic gradients computed on larger mini-batches have smaller variance.
- The mini-batch size $|\mathcal{M}|$ is a hyperparameter that needs to be set.

Stochastic Gradient Descent

- Batch gradient descent moves directly downhill. SGD takes steps in a noisy direction, but moves downhill on average.



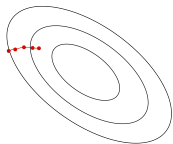
batch gradient descent



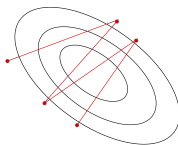
stochastic gradient descent

Learning Rate

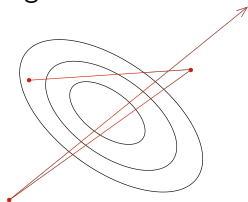
- In gradient descent, the learning rate α is a hyperparameter we need to tune. Here are some things that can go wrong:



α too small:
slow progress



α too large:
oscillations



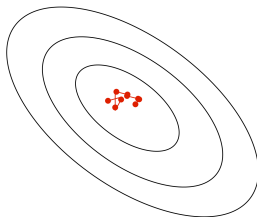
α much too large:
instability

- Good values are typically small. You should do a grid search if you want good performance (i.e. try 0.1, 0.03, 0.01, ...).

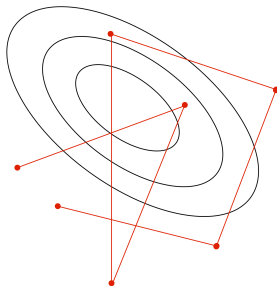
SGD Learning Rate

- In stochastic training, the learning rate also influences the **fluctuations** due to the stochasticity of the gradients.

small learning rate



large learning rate



- Typical strategy:
 - ▶ Use a large learning rate early in training so you can get close to the optimum
 - ▶ Gradually decay the learning rate to reduce the fluctuations

In the last lecture, we have learned the logistic regression for binary classification with $Y \in \{0, 1\}$.

- Estimating the Bayes rule at any observation $\mathbf{x} \in \mathcal{X}$ is equivalent to estimate the conditional probability $\mathbb{P}(Y = 1 \mid X = \mathbf{x})$.
- Logistic regression parametrizes the conditional probability by

$$\mathbb{P}(Y = 1 \mid X = \mathbf{x}) = \frac{e^{\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}}}{1 + e^{\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}}}.$$

- We estimate the coefficients by using MLE which can be solved by (stochastic) gradient descent.

Extension to multi-class classification

When $Y \in \{0, 1, \dots, K\}$ for $K \geq 2$, we need to estimate

$$p_k(\mathbf{x}) := \mathbb{P}(Y = k \mid X = \mathbf{x}), \quad \forall 1 \leq k \leq K.$$

We assume

$$\begin{aligned} p_0(\mathbf{x}) &= \frac{1}{1 + \sum_{k=1}^K e^{\beta_0^{(k)} + \mathbf{x}^\top \boldsymbol{\beta}^{(k)}}}, \\ p_1(\mathbf{x}) &= \frac{e^{\beta_0^{(1)} + \mathbf{x}^\top \boldsymbol{\beta}^{(1)}}}{1 + \sum_{k=1}^K e^{\beta_0^{(k)} + \mathbf{x}^\top \boldsymbol{\beta}^{(k)}}}, \\ &\vdots \\ p_K(\mathbf{x}) &= \frac{e^{\beta_0^{(K)} + \mathbf{x}^\top \boldsymbol{\beta}^{(K)}}}{1 + \sum_{k=1}^K e^{\beta_0^{(k)} + \mathbf{x}^\top \boldsymbol{\beta}^{(k)}}} \end{aligned}$$

Choice of the baseline (which is $Y = 0$) is arbitrary.

Equivalently,

$$\begin{aligned}\log\left(\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}\right) &= \beta_0^{(1)} + \beta_1^{(1)}x_1 + \cdots + \beta_p^{(1)}x_p \\ \log\left(\frac{p_2(\mathbf{x})}{p_0(\mathbf{x})}\right) &= \beta_0^{(2)} + \beta_1^{(2)}x_1 + \cdots + \beta_p^{(2)}x_p \\ &\vdots \\ \log\left(\frac{p_K(\mathbf{x})}{p_0(\mathbf{x})}\right) &= \beta_0^{(K)} + \beta_1^{(K)}x_1 + \cdots + \beta_p^{(K)}x_p\end{aligned}$$

So classification can be done immediately once $\beta^{(k)}$'s are estimated,

How to estimate coefficients?

A naive approach: separate binary logistic regressions

$$\log\left(\frac{p_k(\mathbf{x})}{p_0(\mathbf{x})}\right) = \beta_0^{(k)} + \beta_1^{(k)}x_1 + \dots + \beta_p^{(k)}x_p$$

Split the data into $\{\mathcal{D}_{(1)}^{train}, \dots, \mathcal{D}_{(K)}^{train}\}$ with $\mathcal{D}_{(k)}^{train}$ containing all data with $y \in \{0, k\}$.

1. For each $1 \leq k \leq K$, use $\mathcal{D}_{(k)}^{train}$ to perform binary logistic regression to estimate $\beta^{(k)}$ and estimate

$$\frac{p_k(\mathbf{x})}{p_0(\mathbf{x})}$$

2. Assign class label by comparing

$$1, \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}, \frac{p_2(\mathbf{x})}{p_0(\mathbf{x})}, \dots, \frac{p_K(\mathbf{x})}{p_0(\mathbf{x})}$$

Why naive?

- Estimation of $\beta^{(k)}$
 - ▶ only uses $\mathcal{D}^{train}_{(k)}$, data points in class $\{0, k\}$
 - ▶ ignore all data points in other classes
- The event $\{y_i = k\}$ is **dependent** on all other $\{y_i = k'\}$ for $k' \neq k$. Intuitively, this dependence helps to estimate $\beta^{(k)}$ by pooling data from all classes.
- What should we use instead?

MLE for multi-class logistic regression

For $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$, the log-likelihood of $(\beta^{(1)}, \dots, \beta^{(K)})$ with no intercepts is **proportional to**

$$\begin{aligned} & \sum_{i=1}^n \log \left(\prod_{k=0}^K p_k(\mathbf{x}_i)^{1\{y_i=k\}} \right) \\ &= \sum_{i=1}^n \sum_{k=0}^K 1\{y_i = k\} \log(p_k(\mathbf{x}_i)) \\ &= \sum_{i=1}^n \left[1\{y_i = 0\} \log(p_0(\mathbf{x}_i)) + \sum_{k=1}^K 1\{y_i = k\} \log(p_k(\mathbf{x}_i)) \right] \\ &= \sum_{i=1}^n \left[\sum_{k=1}^K 1\{y_i = k\} \mathbf{x}_i^\top \beta^{(k)} - \sum_{k=0}^K 1\{y_i = k\} \log \left(1 + \sum_{k=1}^K e^{\mathbf{x}_i^\top \beta^{(k)}} \right) \right] \\ &= \sum_{i=1}^n \left[\sum_{k=1}^K 1\{y_i = k\} \mathbf{x}_i^\top \beta^{(k)} - \log \left(1 + \sum_{k=1}^K e^{\mathbf{x}_i^\top \beta^{(k)}} \right) \right] \end{aligned}$$

Gradient of $\ell(\beta^{(k)})$

For any $1 \leq k \leq K$,

$$\begin{aligned}\frac{\partial \ell(\beta^{(1)}, \dots, \beta^{(K)})}{\partial \beta^{(k)}} &= \sum_{i=1}^n \left[1\{y_i = k\} \mathbf{x}_i - \frac{\mathbf{x}_i e^{\mathbf{x}_i^\top \beta^{(k)}}}{1 + \sum_{k=1}^K e^{\mathbf{x}_i^\top \beta^{(k)}}} \right] \\ &= \sum_{i=1}^n \left[1\{y_i = k\} - \frac{e^{\mathbf{x}_i^\top \beta^{(k)}}}{1 + \sum_{k=1}^K e^{\mathbf{x}_i^\top \beta^{(k)}}} \right] \mathbf{x}_i\end{aligned}$$

c.f. the binary case

$$\begin{aligned}\frac{\partial \ell(\beta)}{\partial \beta} &= \sum_{i=1}^n \left[1\{y_i = 1\} - \frac{e^{\mathbf{x}_i^\top \beta}}{1 + e^{\mathbf{x}_i^\top \beta}} \right] \mathbf{x}_i \\ &= \sum_{i=1}^n \left[y_i - \frac{e^{\mathbf{x}_i^\top \beta}}{1 + e^{\mathbf{x}_i^\top \beta}} \right] \mathbf{x}_i.\end{aligned}$$

Therefore, for $1 \leq k \leq K$, we update

$$\hat{\beta}_{(t+1)}^{(k)} = \hat{\beta}_{(t)}^{(k)} + \alpha \sum_{i=1}^n \left[1\{y_i = k\} - \frac{e^{\mathbf{x}_i^\top \hat{\beta}_{(t)}^{(k)}}}{1 + \sum_{k=1}^K e^{\mathbf{x}_i^\top \hat{\beta}_{(t)}^{(k)}}} \right] \mathbf{x}_i.$$

Remark:

- the gradient update uses data points from **all classes**!
- better estimation than the naive approach

An alternative to Logistic Regression

- When the classes are well-separated, the parameter estimates for the logistic regression model are surprisingly unstable¹.
 - ▶ Discriminant analysis does not suffer from this problem.
- When n is small and we know more about the data, such as the distribution of $X \mid Y = k$
 - ▶ Discriminant analysis has better performance than the logistic regression model.
- Logistic Regression sometimes does not handle multi-class classification well
 - ▶ Discriminant analysis is more suitable for **multi-class** classification problems.

¹A paper on this.