

Inferencia variacional

Mauricio A. Álvarez

Curso de entrenamiento ArcelorMittal

Contenido

Introducción

Optimización variacional aplicada a un problema de inferencia

Distribuciones factorizadas

Generalidades

- ❑ Los métodos variacionales tienen sus orígenes en el siglo 18 con el trabajo de Euler, Lagrange y otros, sobre el cálculo de variaciones.
- ❑ En el cálculo estándar se desea encontrar las derivadas de una función.
- ❑ Una función se puede entender como un mapeo que toma una variable de entrada y retorna el valor de la función como su salida.
- ❑ La derivada luego describe cómo cambia el valor de la salida con cambios infinitesimales del valor de la entrada.

Funcional

- Se puede definir un *funcional* como un mapeo que toma una función como su entrada y retorna el valor del funcional como su salida.
- Un ejemplo es la *entropía* $H(p)$, que toma una función de probabilidad $p(x)$ como la entrada y retorna la cantidad

$$H(p) = - \int p(x) \ln p(x) dx,$$

como la salida.

- Se puede introducir el concepto de *derivada funcional*, que expresa cómo cambia el valor del funcional en respuesta a cambios infinitesimales de la función de entrada.

Cálculo de variaciones (I)

- ❑ Las reglas para el cálculo de variaciones son muy parecidas a las reglas del cálculo convencional.
- ❑ Muchos problemas en ciencias e ingeniería, pueden expresarse en términos de un problema de optimización en el que la cantidad que se desea optimizar es un *funcional*.
- ❑ La solución se obtiene explorando todas las posibles funciones de entrada que maximizan o minimizan el funcional.
- ❑ Los métodos variacionales se emplean para encontrar soluciones aproximadas a este tipo de problemas de optimización.

Cálculo de variaciones (II)

- ❑ Esto se realiza restringiendo la clase de funciones sobre las cuales se realiza la optimización.
- ❑ Por ejemplo, considerando sólo funciones cuadráticas o considerando funciones compuestas por funciones base fijas.
- ❑ En inferencia probabilística, la restricción puede tomar la forma de una factorización.

Contenido

Introducción

Optimización variacional aplicada a un problema de inferencia

Distribuciones factorizadas

Modelo Bayesiano

- Supongamos que se tiene un modelo Bayesiano completo en el que a todos los parámetros se les ha asignado distribuciones prior.
- El modelo podría tener tanto variables latentes como parámetros, conjuntamente denotados como \mathbf{Z} .
- Similarmente, se denotan las variables observadas como \mathbf{X} .
- El modelo probabilístico especifica la función de distribución conjunta $p(\mathbf{X}, \mathbf{Z})$ y el objetivo es encontrar una aproximación a la distribución posterior $p(\mathbf{Z}|\mathbf{X})$ como a la evidencia del modelo $p(\mathbf{X})$.

Probabilidad marginal logarítmica

- La probabilidad marginal logarítmica $\ln p(\mathbf{X})$ se puede escribir como

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q\|p),$$

donde

$$\mathcal{L}(q) = \int_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$
$$\text{KL}(q\|p) = - \int_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} \right\} d\mathbf{Z}.$$

donde $q(\mathbf{Z})$ es la distribución desconocida.

- Se puede maximizar el límite inferior $\mathcal{L}(q)$ optimizando con respecto a la distribución $q(\mathbf{Z})$, que es equivalente a minimizar la divergencia de Kullback-Leibler (KL).

¿Cómo escoger $q(\mathbf{Z})$? (I)

- Si se permite cualquier forma para $q(\mathbf{Z})$, el máximo del límite inferior ocurre cuando la divergencia KL se hace cero, que a su vez ocurre cuando $q(\mathbf{Z})$ iguala a la distribución posterior $p(\mathbf{Z}|\mathbf{X})$.
- En la práctica, se asume que trabajar con la verdadera distribución posterior es intratable.
- En su lugar se considera una familia de distribuciones restringidas $q(\mathbf{Z})$, para las cuales se minimice la divergencia KL.

¿Cómo escoger $q(\mathbf{Z})$? (II)

- Una forma de restringir la familia de distribuciones es usar una distribución paramétrica $q(\mathbf{Z}|\omega)$, gobernada por un conjunto de parámetros ω .
- El límite inferior $\mathcal{L}(q)$ se vuelve entonces una función de ω , y se pueden explotar técnicas de optimización no lineal para determinar los valores óptimos de los parámetros.

Contenido

Introducción

Optimización variacional aplicada a un problema de inferencia

Distribuciones factorizadas

Mean field (I)

- Supongamos que los elementos de \mathbf{Z} se dividen en grupos que no se traslapan, que se denotan como $\mathbf{Z}_i, i = 1, \dots, M$.
- Luego se asume que la distribución $q(\mathbf{Z})$ se factoriza con respecto a estos grupos, tal que

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i).$$

- Nótese que esta es la única restricción que se hace sobre $q(\mathbf{Z})$.
- Esta forma factorizada de inferencia variacional corresponde a una aproximación desarrollada en física conocida como “mean field theory”.

Mean field (II)

- Entre todas las distribuciones $q(\mathbf{Z})$ que tienen la forma factorizada anterior, se busca aquella distribución para la cual el límite inferior $\mathcal{L}(q)$ sea el mayor.
- Se desea realizar una optimización variacional de $\mathcal{L}(q)$, con respecto a todas las distribuciones $q_i(\mathbf{Z}_i)$, que se puede realizar optimizando con respecto a cada uno de los factores a la vez.

$\mathcal{L}(q)$ en función de $q_i(\mathbf{Z}_i)$

El límite inferior $\mathcal{L}(q)$ se puede escribir como

$$\begin{aligned}\mathcal{L}(q) &= \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z} \\ &= \int \left(\prod_{i=1}^M q_i(\mathbf{Z}_i) \right) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{\prod_{i=1}^M q_i(\mathbf{Z}_i)} \right\} d\mathbf{Z} \\ &= \int \prod_{\forall i} q_i \{ \ln p(\mathbf{X}, \mathbf{Z}) \} d\mathbf{Z} - \int \prod_{\forall i} q_i \sum_{\forall i} \ln q_i d\mathbf{Z}\end{aligned}$$

Distribuciones óptimas $q_j^*(\mathbf{Z}_j)$

- Se puede demostrar que la solución óptima $q_j^*(\mathbf{Z}_j)$ está dada como

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{const},$$

donde $\mathbb{E}_{i \neq j}$ se toma con respecto a todos los factores q_i con $i \neq j$.

- La constante aditiva const se selecciona de forma tal que normalice la distribución $q_j^*(\mathbf{Z}_j)$.
- Tomando la exponencial en ambos lados, en la expresión anterior, se tiene

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j}.$$

- En la práctica, se prefiere trabajar con la expresión del $\ln q_j^*(\mathbf{Z}_j)$ y después normalizar.

Procedimiento

- El valor óptimo de $q_j^*(\mathbf{Z}_j)$ depende de los valores esperados calculados con respecto a los otros factores $q_i(\mathbf{Z}_i)$, para $i \neq j$.
- La solución se obtiene inicializando primero todos los factores $q_i(\mathbf{Z}_i)$, y luego iterando a través de cada factor.
- Cada factor se actualiza usando la expresión obtenida para $q_i^*(\mathbf{Z}_i)$.
- Para actualizar cualquier expresión, se usan los factores que ya se hayan actualizado hasta ese momento.