

Прикладная статистика в машинном обучении  
Семинар 9  
Байесовский взгляд на подбор моделей

И. К. Козлов  
(Мехмат МГУ)

2022

Сразу съедем лягушку

Настало время съесть лягушку!



Обсудим достаточные статистики экспоненциального класса!

# Достаточные статистики

## 1 Достаточные статистики

## 2 Экспоненциальное семейство распределений

## 3 Дифференцирование

- Дифференциал
- Свойства дифференциала
- Примеры вычисления
- Гессиан
- Памятка

### Теорема факторизации Фишера

$T(x)$  — достаточная статистика  $\Leftrightarrow$  существует разложение

$$f(x | \theta) = h(x) q(\theta, T(x)).$$

Докажем эту теорему в **дискретном случае**. В общем случае — см. Wikipedia:

[https://en.wikipedia.org/wiki/Sufficient\\_statistic](https://en.wikipedia.org/wiki/Sufficient_statistic)

### Доказательство теоремы факторизации.

( $\Rightarrow$ ). Пусть  $S(x)$  — достаточная, т.е.  $\mathbb{P}(X \mid T = t)$  не зависит от  $\theta$ .

Обозначим  $t = T(x)$ . Тогда

$$f(x \mid \theta) = \mathbb{P}(X = x) = \underbrace{\mathbb{P}(X = x \mid T(X) = t)}_{\text{не зависит от } \theta} \underbrace{\mathbb{P}(T = t)}_{\text{функция от } \theta \text{ и } T}.$$

Это и есть нужное разложение:

$$f(x \mid \theta) = h(x) q(\theta, T(x)).$$

## Достаточные статистики

( $\Leftarrow$ ). (В дискретном случае). Пусть

$$f(x \mid \theta) = h(x) q(\theta, T(x)). \quad (1)$$

По определению условной вероятности

$$\mathbb{P}(X = x \mid T = t) = \frac{\mathbb{P}(X = x, T = t)}{\mathbb{P}(T = t)} = \frac{\mathbb{P}(X = x)}{\mathbb{P}(T = t)}.$$

## Достаточные статистики

( $\Leftarrow$ ). (В дискретном случае). Пусть

$$f(x \mid \theta) = h(x) q(\theta, T(x)). \quad (1)$$

По определению условной вероятности

$$\mathbb{P}(X = x \mid T = t) = \frac{\mathbb{P}(X = x, T = t)}{\mathbb{P}(T = t)} = \frac{\mathbb{P}(X = x)}{\mathbb{P}(T = t)}.$$

Числитель равен  $f(x; \theta)$ , а знаменатель — сумма вероятностей событий

$$\mathbb{P}(T = t) = \sum_{x': T(x')=t} f(x'; \theta)$$

Подставляем (1):

$$\mathbb{P}(X = x \mid T = t) = \frac{f(x; \theta)}{\sum_{x': T(x')=t} f(x'; \theta)} = \frac{q(\theta, t) h(x)}{q(\theta, t) \sum_{x': T(x')=t} h(x')}$$

## Достаточные статистики

( $\Leftarrow$ ). (В дискретном случае). Пусть

$$f(x \mid \theta) = h(x) q(\theta, T(x)). \quad (1)$$

По определению условной вероятности

$$\mathbb{P}(X = x \mid T = t) = \frac{\mathbb{P}(X = x, T = t)}{\mathbb{P}(T = t)} = \frac{\mathbb{P}(X = x)}{\mathbb{P}(T = t)}.$$

Числитель равен  $f(x; \theta)$ , а знаменатель — сумма вероятностей событий

$$\mathbb{P}(T = t) = \sum_{x': T(x')=t} f(x'; \theta)$$

Подставляем (1):

$$\mathbb{P}(X = x \mid T = t) = \frac{f(x; \theta)}{\sum_{x': T(x')=t} f(x'; \theta)} = \frac{q(\theta, t) h(x)}{q(\theta, t) \sum_{x': T(x')=t} h(x')}$$

$q(\theta, t)$  сокращается. Что и требовалось —  $\mathbb{P}(X = x \mid T(X) = t)$  не зависит от  $\theta$ .



# Экспоненциальное семейство распределений

1 Достаточные статистики

2 Экспоненциальное семейство распределений

3 Дифференцирование

- Дифференциал
- Свойства дифференциала
- Примеры вычисления
- Гессиан
- Памятка

### Экспоненциальный класс

Семейство распределений, чья плотность может быть представлена в виде:

$$f(x|\theta) = \frac{h(x)}{g(\theta)} \exp\left(\theta^T u(x)\right),$$

- $h(x) \geq 0$ ,
- $\theta^T u(x)$  означает  $\theta_1 u_1(x) + \dots + \theta_k u_k(x)$ ,
- Чтобы интеграл от плотности равнялся единице:

$$g(\theta) = \int h(x) \exp\left(\theta^T u(x)\right) dx.$$

## Экспоненциальное семейство

Как мы помним,  $u_j(x)$  — достаточные статистики:

$$f(x|\theta) = \frac{h(x)}{g(\theta)} \exp\left(\theta^T u(x)\right).$$

## Экспоненциальное семейство

Как мы помним,  $u_j(x)$  — достаточные статистики:

$$f(x|\theta) = \frac{h(x)}{g(\theta)} \exp\left(\theta^T u(x)\right).$$

Их моменты выражаются через производную  $g(\theta)$ :

- Матожидание

$$\mathbb{E}(u_j) = \frac{\partial \ln g(\theta)}{\partial \theta_j}.$$

- Матрица ковариации

$$\text{Cov}(u_i, u_j) = \frac{\partial^2 \ln g(\theta)}{\partial \theta_i \partial \theta_j}.$$

## Экспоненциальное семейство

**Доказательство.** Докажем для матожидания, для ковариации — аналогично.

- Q: Если  $f(x|\theta) = \frac{h(x)}{g(\theta)} \exp(\theta^T u(x))$ , то чему равна  $g(\theta)$ ?

## Экспоненциальное семейство

**Доказательство.** Докажем для матожидания, для ковариации — аналогично.

- **Q:** Если  $f(x|\theta) = \frac{h(x)}{g(\theta)} \exp(\theta^T u(x))$ , то чему равна  $g(\theta)$ ?

- **A:** Это нормировочная константа:

$$g(\theta) = \int_{-\infty}^{\infty} h(x) \exp(\theta^T u(x)) dx.$$

- Ну конечно, интеграл можно продифференцировать по параметру:

$$\frac{\partial \ln g(\theta)}{\partial \theta_j} = \frac{1}{g(\theta)} \frac{\partial g(\theta)}{\partial \theta_j} = \int_{-\infty}^{\infty} u_j(x) \frac{h(x)}{g(\theta)} \exp(\theta^T u(x)) dx = \mathbb{E}(u_j(x)).$$

Вычисления упрощаются

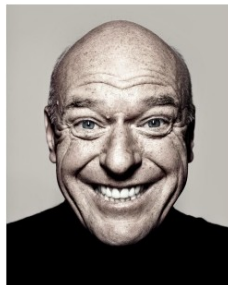
$$\mathbb{E}(u_j)$$

**Интеграл  
считать**



$$\frac{\partial \log g(\theta)}{\partial \theta_j}$$

**Производную  
взять**



# Дифференцирование

1 Достаточные статистики

2 Экспоненциальное семейство распределений

3 Дифференцирование

- Дифференциал
- Свойства дифференциала
- Примеры вычисления
- Гессиан
- Памятка



# Дифференциал

- 1 Достаточные статистики
- 2 Экспоненциальное семейство распределений
- 3 Дифференцирование
  - Дифференциал
  - Свойства дифференциала
  - Примеры вычисления
  - Гессиан
  - Памятка

## Дифференцирование

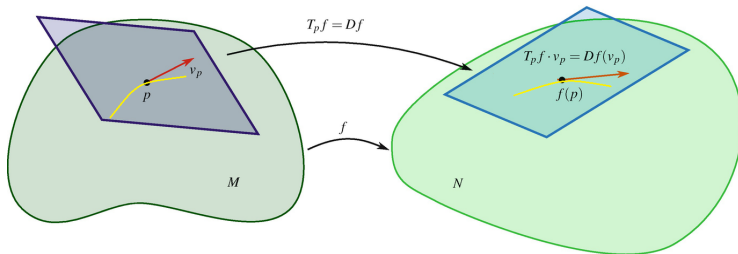
Наши мучения не заканчиваются)  
Давайте вспомним — как **дифференцировать функции**.



# Дифференцирование

Пусть  $f : M^n \rightarrow N^m$  — гладкое отображение. Тогда в каждой точке  $p \in M^n$  определён дифференциал отображения

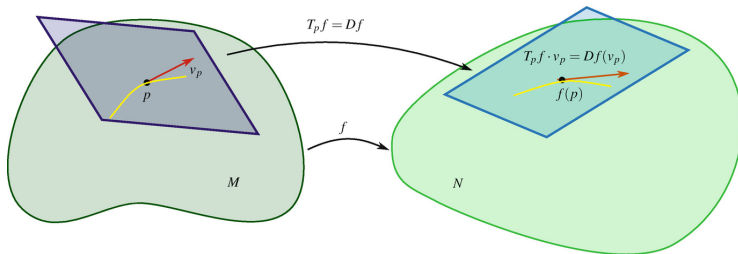
$$d_p f : T_p M^n \rightarrow T_{f(p)} N^m.$$



## Дифференцирование

Пусть  $f : M^n \rightarrow N^m$  — гладкое отображение. Тогда в каждой точке  $p \in M^n$  определён дифференциал отображения

$$d_p f : T_p M^n \rightarrow T_{f(p)} N^m.$$



Жёлтая кривая  $\gamma_M(t)$  на  $M$  переходит в жёлтую кривую  $\gamma_N(t)$  на  $N$ .

Дифференциал  $d_p f$  переводит вектор скорости  $\frac{d\gamma_M}{dt}$  в вектор скорости  $\frac{d\gamma_N}{dt}$ .

# Дифференцирование

Мы будем считать, что  $M^n$  и  $N^m$  — области в  $\mathbb{R}^n$  и  $\mathbb{R}^m$  соответственно.

Поэтому дифференциал будет линейным отображением

$$d_p f : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

Будем писать  $d_p f = df$  для краткости.

Как правильно вычислять в криволинейных координатах — учит  
дифференциальная геометрия.

# Дифференцирование

Рассмотрим примеры:

- **Функция.**  $f(x) : \mathbb{R} \rightarrow \mathbb{R}$ . Какой вид имеет дифференциал?

# Дифференцирование

Рассмотрим примеры:

- **Функция.**  $f(x) : \mathbb{R} \rightarrow \mathbb{R}$ . Какой вид имеет дифференциал?

$$df(x) = f'(x)dx$$

- Функцию от вектора  $f(x^1, \dots, x^n) : \mathbb{R}^n \rightarrow \mathbb{R}$ . Дифференциал?



- **Функцию от вектора**  $f(x^1, \dots, x^n) : \mathbb{R}^n \rightarrow \mathbb{R}$ . Дифференциал?

$$df(x) = \sum_{i=1}^n \frac{\partial f}{\partial x^i} dx^i.$$

Вектор частных производных<sup>1</sup> — **градиент функции**:

$$\nabla f = \left( \frac{\partial f}{\partial x^1}, \dots, \frac{\partial f}{\partial x^n} \right).$$

Поэтому дифференциал — формальное скалярное произведение

$$df(x) = \langle \nabla f, dx \rangle$$

---

<sup>1</sup> в евклидовых координатах

- Далее — дифференциал **вектор-функции от вектора**

$$y^i = f^i(x^1, \dots, x^n), \quad i = 1, \dots, m.$$

# Дифференцирование

- Далее — дифференциал **вектор-функции от вектора**

$$y^i = f^i(x^1, \dots, x^n), \quad i = 1, \dots, m.$$

Дифференциал

$$df(x) = J_f(x)dx,$$

где **матрица Якоби**

$$J_f = \begin{pmatrix} \frac{\partial f^1}{\partial x^1} & \dots & \frac{\partial f^1}{\partial x^n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f^m}{\partial x^1} & \dots & \frac{\partial f^m}{\partial x^n} \end{pmatrix}.$$

- **Функция от матрицы**  $f(X)$ , где  $X$  —  $m \times n$  матрица.

**Q:** Какой размер у матрицы дифференциала?

# Дифференцирование

- **Функция от матрицы**  $f(X)$ , где  $X$  —  $m \times n$  матрица.

**Q:** Какой размер у матрицы дифференциала?

**A:**  $m \times n$  как у  $X$ . Пусть  $X^{ij}$  — компоненты  $X$ . Дифференциал:

$$df(X) = \sum_{i,j} \frac{\partial f}{\partial X^{ij}} dX^{ij} = \langle \nabla f, dX \rangle$$

## Утверждение

Для любых  $m \times n$  матриц  $A, B$

$$\langle A, B \rangle = \text{tr}(A^T B).$$

## Утверждение

Для любых  $m \times n$  матриц  $A, B$

$$\langle A, B \rangle = \text{tr}(A^T B).$$

**Доказательство.** Обозначим  $A = (a_{ij})$ ,  $B = (b_{ij})$  и  $C = A^T B = (c_{ij})$ . Тогда

$$c_{ij} = \sum_{k=1}^m a_{ki} b_{kj}.$$

## Утверждение

Для любых  $m \times n$  матриц  $A, B$

$$\langle A, B \rangle = \operatorname{tr}(A^T B).$$

**Доказательство.** Обозначим  $A = (a_{ij})$ ,  $B = (b_{ij})$  и  $C = A^T B = (c_{ij})$ . Тогда

$$c_{ij} = \sum_{k=1}^m a_{ki} b_{kj}.$$

Поэтому

$$\operatorname{tr}(A^T B) = \sum_{i=1}^n \sum_{k=1}^m a_{ki} b_{ki} = \langle A, B \rangle.$$



Итого:

Вход \ Выход	Скаляр	Вектор	Матрица
Скаляр	$df(x) = f'(x)dx$ ( $f'(x)$ : скаляр; $dx$ : скаляр)	—	—
Вектор	$df(x) = \langle \nabla f(x), dx \rangle$ ( $\nabla f(x)$ : вектор; $dx$ : вектор)	$df(x) = J_f(x)dx$ ( $J_f(x)$ : матрица; $dx$ : вектор)	—
Матрица	$df(X) = \langle \nabla f(X), dX \rangle$ ( $\nabla f(X)$ : матрица; $dX$ : матрица)	—	—

# Свойства дифференциала

1 Достаточные статистики

2 Экспоненциальное семейство распределений

3 Дифференцирование

- Дифференциал
- Свойства дифференциала
- Примеры вычисления
- Гессиан
- Памятка

## Свойства дифференциала.

❶ Это линейный функционал:

$$d(X + Y) = dX + dY, \quad d(\alpha X) = \alpha d(X).$$

Если  $A, B$  — постоянные матрицы, то

$$d(AXB) = Ad(X)B.$$

## ② Правило Лейбница

$$d(fg) = (df)g + f(dg)$$

С обычного произведения оно продолжается на матричные произведения:

$$d(XY) = (dX)Y + X(dY)$$

и скалярные произведения

$$d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle.$$

## ③ Производная сложной функции.

Дифференциал композиции отображений

$$M \xrightarrow{f} N \xrightarrow{g} P$$

будет композицией дифференциалов

$$d_x f(g(x)) = d_{f(x)} g \circ d_x f(x).$$

## 3 Производная сложной функции.

Дифференциал композиции отображений

$$M \xrightarrow{f} N \xrightarrow{g} P$$

будет композицией дифференциалов

$$d_x f(g(x)) = d_{f(x)} g \circ d_x f(x).$$

В координатах

$$\frac{\partial f(g(x))}{\partial x^i} = \sum_j \frac{\partial f}{\partial g^j} \frac{\partial g^j}{\partial x^i}.$$

# Примеры вычисления

1 Достаточные статистики

2 Экспоненциальное семейство распределений

3 Дифференцирование

- Дифференциал
- Свойства дифференциала
- **Примеры вычисления**
- Гессиан
- Памятка

# Дифференцирование

Q: Найти дифференциал

$$d(X^{-1}) = ?$$



# Дифференцирование

**Q:** Найти дифференциал

$$d(X^{-1}) = ?$$

**A:** Дифференцируем тождество

$$X^{-1}X = E.$$

По правилу Лейбница

$$d(X^{-1})X + X^{-1}dX = 0,$$

откуда

$$d(X^{-1}) = -X^{-1}dX X^{-1}.$$

# Дифференцирование

**Q:** Найти градиент

$$\nabla (\text{tr } A X B X^{-1}) = ?$$

## Дифференцирование

**Q:** Найти градиент

$$\nabla (\operatorname{tr} A X B X^{-1}) = ?$$

**A:** Найдём его из свойства

$$\operatorname{tr} \left[ (\nabla f)^T dX \right].$$

След — линейный оператор, поэтому

$$d(\operatorname{tr} X) = \operatorname{tr}(dX).$$

Далее используем правило Лейбница

$$d(\operatorname{tr} A X B X^{-1}) = \operatorname{tr}(A d(X) B X^{-1}) + \operatorname{tr}(A X B d(X^{-1})).$$

## Дифференцирование

Подставляя  $d(X^{-1})$  получаем

$$d(\operatorname{tr} AXBX^{-1}) = \operatorname{tr}(Ad(X)BX^{-1}) - \operatorname{tr}(AXBX^{-1}d(X)X^{-1}).$$

## Дифференцирование

Подставляя  $d(X^{-1})$  получаем

$$d(\operatorname{tr} AXBX^{-1}) = \operatorname{tr}(Ad(X)BX^{-1}) - \operatorname{tr}(AXBX^{-1}d(X)X^{-1}).$$

Приводим выражение к виду

$$\operatorname{tr}\left[(\nabla f)^T dX\right],$$

используя свойства следа

$$\operatorname{tr} A = \operatorname{tr} A^T, \quad \operatorname{tr} AB = \operatorname{tr} BA.$$

## Дифференцирование

Подставляя  $d(X^{-1})$  получаем

$$d(\operatorname{tr} AXBX^{-1}) = \operatorname{tr}(Ad(X)BX^{-1}) - \operatorname{tr}(AXBX^{-1}d(X)X^{-1}).$$

Приводим выражение к виду

$$\operatorname{tr}[(\nabla f)^T dX],$$

используя свойства следа

$$\operatorname{tr} A = \operatorname{tr} A^T, \quad \operatorname{tr} AB = \operatorname{tr} BA.$$

Получаем

$$d(\operatorname{tr} AXBX^{-1}) = \operatorname{tr}\left[\left(A^T X^{-T} B^T A - X^{-T} B^T X^T A^T X^{-T}\right)^T dX\right]$$

Ответ:

$$\nabla f = A^T X^{-T} B^T A - X^{-T} B^T X^T A^T X^{-T}.$$

# Дифференцирование

Q: Найти дифференциал

$$d \ln (\det X) = ?$$

# Дифференцирование

**Q:** Найти дифференциал

$$d \ln (\det X) = ?$$

**A:** Производная сложной функции:

$$d \ln (\det X) = \frac{d (\det X)}{\det X}.$$

Далее вспомним пару фактов о матрицах.



# Дифференцирование

Рассмотрим произвольную матрицу  $A = (a_{ij})$ .

Обозначим  $A_{ij}$  — алгебраические дополнения к  $a_{ij}$ .

# Дифференцирование

Рассмотрим произвольную матрицу  $A = (a_{ij})$ .

Обозначим  $A_{ij}$  — алгебраические дополнения к  $a_{ij}$ .

- $\det A$  — многочлен от  $a_{ij}$ ;
- Разложения определителя по строке/столбцу:

$$\det A = \sum_{j=1}^n a_{ij} A_{ij} = \sum_{i=1}^n a_{ij} A_{ij}.$$

# Дифференцирование

Рассмотрим произвольную матрицу  $A = (a_{ij})$ .

Обозначим  $A_{ij}$  — алгебраические дополнения к  $a_{ij}$ .

- $\det A$  — многочлен от  $a_{ij}$ ;
- Разложения определителя по строке/столбцу:

$$\det A = \sum_{j=1}^n a_{ij} A_{ij} = \sum_{i=1}^n a_{ij} A_{ij}.$$

- ❶ частные производные определителя — алгебраические дополнения

$$\frac{\partial \det A}{\partial a_{ij}} = A_{ij}, \quad (2)$$

# Дифференцирование

- ❶ частные производные определителя — алгебраические дополнения

$$\frac{\partial \det A}{\partial a_{ij}} = A_{ij}, \quad (2)$$

- ❷ элементы обратной матрицы — это алгебраические дополнения  $A_{ji}$ , делённые на определитель:

$$A^{-1} = \frac{1}{\det(A)} \operatorname{adj}(A), \quad \text{где} \quad \operatorname{adj}(A)_{ij} = A_{ji}.$$

## Дифференцирование

Ответ:

$$d \ln (\det X) = \langle X^{-T}, dX \rangle.$$

- 1 Достаточные статистики
- 2 Экспоненциальное семейство распределений
- 3 Дифференцирование
  - Дифференциал
  - Свойства дифференциала
  - Примеры вычисления
  - Гессиан
  - Памятка

# Гессиан

Матрицу вторых производных называют **гессианом** и обозначают  $\mathbf{H}_f$  или  $\nabla^2 f$ :

$$\nabla^2 f = \left( \frac{\partial^2 f}{\partial x_i \partial x_j} \right).$$



## Формальное вычисление гессиана.

- Берём дифференциал

$$df(x) = \langle \nabla f(x), dx \rangle.$$

- Меняем  $dx \rightarrow dx_1$ .

- Ещё раз берём дифференциал

$$d\langle \nabla f(x), dx_1 \rangle = \langle d(\nabla f(x)), dx_1 \rangle = \langle \nabla^2 f(x) dx_2, dx_1 \rangle.$$

## Гессиан

**Q:** Найти гессиан квадратичной функции

$$f(x) = \frac{1}{2} \langle x, Ax \rangle + b^T x + c, \quad x \in \mathbb{R}^n.$$

## Гессиан

**Q:** Найти гессиан квадратичной функции

$$f(x) = \frac{1}{2} \langle x, Ax \rangle + b^T x + c, \quad x \in \mathbb{R}^n.$$

**A:** Находим дифференциал

$$df(x) = \langle Ax + b, dx \rangle.$$

И ещё раз берём дифференциал

$$d(Ax + b) = A.$$

**Ответ:**  $\nabla^2 f = A$ .

- 1 Достаточные статистики
- 2 Экспоненциальное семейство распределений
- 3 Дифференцирование
  - Дифференциал
  - Свойства дифференциала
  - Примеры вычисления
  - Гессиан
  - **Памятка**

## Правила преобразования

$$dA = 0$$

$$d(\alpha X) = \alpha(dX)$$

$$d(AXB) = A(dX)B$$

$$d(X + Y) = dX + dY$$

$$d(X^T) = (dX)^T$$

$$d(XY) = (dX)Y + X(dY)$$

$$d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$$

$$d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$$

## Таблица стандартных производных

$$d\langle A, X \rangle = \langle A, dX \rangle$$

$$d\langle Ax, x \rangle = \langle (A + A^T)x, dx \rangle$$

$$d(\operatorname{tr} X) = \operatorname{tr}(dX)$$

$$d(\operatorname{Det}(X)) = \operatorname{Det}(X)\langle X^{-T}, dX \rangle$$

$$d(X^{-1}) = -X^{-1}(dX)X^{-1}$$