

# Машинное обучение и нейросетевые модели

## *Лекция 8. Вероятностные модели и их применение*

Лектор: Кравченя Павел Дмитриевич

Волгоград 2025

---

## План лекции

1. Недостатки линейных моделей. Переход к обобщенным линейным моделям.
  2. Функция связи. Экспоненциальный класс распределений.
  3. Лемма о градиенте кумулятивной функции.
  4. Подкласс экспоненциальных распределений. Каноническая функция связи.
  5. Байесовские линейная, логистическая и пуассоновская регрессии.
  6. Понятие моделей смесей. Структура моделей смесей.
  7. Смеси гауссовских и категориальных моделей. LDA.
  8. Байесовские нейронные сети. Стохастические нейронные сети.
  9. Непараметрические модели. Гауссовы процессы. Байесовская оптимизация.
  10. Реализация байесовской нейронной сети во фреймворке Pyro.
-

- Одним из широко используемых классов моделей являются линейные, в которых *матожидание* целевой переменной определяется линейной комбинацией признаков:

$$\mathbb{E}(y \mid \mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle.$$

- Однако, у линейных моделей есть существенные недостатки:
  1. Линейные модели *не позволяют* учесть нелинейную связь между признаками и целевой переменной.
  2. Дисперсия линейных моделей, как правило, является *константной* и *не зависящей* от признаков.
  3. Линейные модели *не позволяют* обрабатывать дискретную целевую переменную.

- Попробуем обобщить *линейные модели*, введя некоторую монотонную функцию  $g(t)$ , связывающую *матожидание* целевой переменной и *линейную комбинацию* признаков в модели:

$$g(\mathbb{E}(y | \mathbf{x})) = \langle \mathbf{w}, \mathbf{x} \rangle.$$

- Данная функция  $g(t)$  называется функцией связи (*link function*).
- Таким образом, функция связи позволяет линеаризовать матожидание целевой переменной по признакам.
- Таким образом, для определения конкретной модели требуются:
  - Параметризованное семейство распределений:  $p_{\theta}(y | \langle \mathbf{x}, \mathbf{w} \rangle)$ .
  - Функцию связи  $g(t)$ .
- Определенная таким образом модель называется обобщенной линейной.



- При построении обобщенной линейной модели можно выбрать любой класс распределений  $p_{\theta}(y | \langle \mathbf{x}, \mathbf{w} \rangle)$  и любую монотонную функцию связи  $g(t)$ , получив при этом некоторую вероятностную модель.
- Однако, часто предполагают, что класс распределений принадлежит одному из достаточно простых семейств экспоненциального класса, который в общем виде выражается следующим образом:

$$p(y | \boldsymbol{\theta}) = \frac{1}{h(\boldsymbol{\theta})} f(y) \cdot \exp(\boldsymbol{\theta}^T \mathbf{u}(y)).$$

- Здесь  $\boldsymbol{\theta}$  – вектор вещественнозначных параметров (значениями которых различаются распределения из семейства),  $h > 0$ ,  $f > 0$ ,  $\mathbf{u}$  – некоторая вектор-функция.

- Для экспоненциального семейства:

$$p(y | \boldsymbol{\theta}) = \frac{1}{h(\boldsymbol{\theta})} f(y) \cdot \exp(\boldsymbol{\theta}^T \mathbf{u}(y)).$$

- Интегрируя обе части равенства по  $y$  и применяя условие нормировки, получаем:

$$h(\boldsymbol{\theta}) = \int f(y) \cdot \exp(\boldsymbol{\theta}^T \mathbf{u}(y)) dy.$$

- Вычислим градиент от логарифма  $h(\boldsymbol{\theta})$ :

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \log h(\boldsymbol{\theta}) &= \frac{\nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta})}{h(\boldsymbol{\theta})} = \frac{\nabla_{\boldsymbol{\theta}} \int f(y) \cdot \exp(\boldsymbol{\theta}^T \mathbf{u}(y)) dy}{h(\boldsymbol{\theta})} = \frac{\int f(y) \nabla_{\boldsymbol{\theta}} \exp(\boldsymbol{\theta}^T \mathbf{u}(y)) dy}{h(\boldsymbol{\theta})} = \\ &= \frac{\int f(y) \mathbf{u}(y) \exp(\boldsymbol{\theta}^T \mathbf{u}(y)) dy}{h(\boldsymbol{\theta})} = \int \mathbf{u}(y) \cdot \frac{1}{h(\boldsymbol{\theta})} f(y) \exp(\boldsymbol{\theta}^T \mathbf{u}(y)) dy = \mathbb{E}_{p(y|\boldsymbol{\theta})}[\mathbf{u}(y)]. \end{aligned}$$

- В контексте *обобщенных линейных моделей* обычно рассматривают подкласс экспоненциального класса, состоящий из семейств, представимых в виде:

$$p(y | \theta, \phi) = \exp \left( \frac{y\theta - a(\theta)}{\phi} + b(y, \phi) \right).$$

где  $\theta$  и  $\phi$  – скалярные параметры, причем  $\phi$  – фиксированный параметр, а значения  $\theta$  параметризуют распределения из семейства. Величины  $a(\theta)$  и  $b(y, \phi)$  – некоторые скалярные функции.

- Данное выражение можно переписать в виде, больше похожим на канонический вид экспоненциального класса распределений:

$$p(y | \theta, \phi) = \frac{1}{\exp \left( \frac{a(\theta)}{\phi} \right)} \exp[b(y, \phi)] \cdot \exp \left( \frac{y\theta}{\phi} \right).$$

- Поскольку  $\phi$  является фиксированным параметром (константой), имеем:

$$p(y | \theta, \phi) = \frac{1}{h(\theta)} f(y) \cdot \exp(\boldsymbol{\theta}^T \mathbf{u}(y)), \quad h(\theta) = \exp\left(\frac{a(\theta)}{\phi}\right),$$

- Видно, что функция  $\mathbf{u}(y)$  содержит только одну компоненту  $u_1(y)$ , причём:

$$u_1(y) = \frac{y}{\phi}.$$

- В соответствии с доказанной леммой получаем:

$$\mathbb{E}[u_1(y)] = \nabla_{\theta} \log h(\theta),$$

$$\frac{\mathbb{E}[y]}{\phi} = \frac{\partial}{\partial \theta} \left( \frac{a(\theta)}{\phi} \right) = \frac{1}{\phi} \frac{\partial}{\partial \theta} (a(\theta)).$$

$$\mathbb{E}[y] = \frac{\partial}{\partial \theta} (a(\theta)) = a'(\theta).$$



- Как будет выглядеть выражение для  $p(y | \mathbf{x})$ ?
- Единственное, что может *зависеть* от  $\mathbf{x}$  в выражении *экспоненциального класса*, – это параметр  $\theta$ . Положим  $\theta = \langle \mathbf{w}, \mathbf{x} \rangle$ . Тогда:

$$\mathbb{E}[y | \mathbf{x}] = a'(\langle \mathbf{w}, \mathbf{x} \rangle).$$

- Но для функции связи справедливо соотношение:

$$g(\mathbb{E}(y | \mathbf{x})) = \langle \mathbf{w}, \mathbf{x} \rangle.$$

- Отсюда можно однозначно определить функцию связи:

$$g = (a')^{-1}.$$

- *Функция связи*, введенная таким образом, называется канонической.
- Поэтому, чтобы определить функцию связи, нужно вычислить  $a'(\langle \mathbf{w}, \mathbf{x} \rangle)$ .

- Рассмотрим байесовскую линейную регрессию:

$$p(y | \mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[ -\frac{(y - \langle \mathbf{w}, \mathbf{x} \rangle)^2}{2\sigma^2} \right].$$

- Вычислим функцию связи для данного семейства распределений. Оно принадлежит к экспоненциальному классу. Попробуем представить его в виде:

$$p(y | \mu, \sigma^2) = \exp \left( \frac{y\mu - a(\mu)}{\sigma^2} + b(y, \sigma^2) \right).$$

$$\begin{aligned} p(y | \mu, \sigma^2) &= \exp \left[ -\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(y - \mu)^2}{2\sigma^2} \right] = \exp \left[ -\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{y^2}{2\sigma^2} + \frac{y\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} \right] = \\ &= \exp \left[ \frac{y\mu - (\mu^2/2)}{\sigma^2} - \left( \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{y^2}{2\sigma^2} \right) \right] \Rightarrow a(\mu) = \frac{\mu^2}{2}. \end{aligned}$$

- Тогда  $a'(\mu) = \mu$ ,  $\mathbb{E}(y | \mathbf{w}, \mathbf{x}) = g^{-1}(\langle \mathbf{w}, \mathbf{x} \rangle) = a'(\langle \mathbf{w}, \mathbf{x} \rangle) = \langle \mathbf{w}, \mathbf{x} \rangle$ .

- Пусть  $\mu = \mathbb{E}[\langle \mathbf{x}, \mathbf{w} \rangle]$ . Рассмотрим байесовскую логистическую регрессию:

$$p(y | \mu) = \mu^y (1 - \mu)^{1-y}.$$

- Здесь используется семейство распределений Бернулли, которое тоже принадлежит к экспоненциальному классу.

$$\begin{aligned} p(y | \mu) &= \mu^y (1 - \mu)^{1-y} = \exp[\log(\mu^y \cdot (1 - \mu)^{1-y})] = \exp \left[ \log \left( \frac{\mu^y}{(1 - \mu)^{y-1}} \right) \right] = \\ &= \exp \left[ \log \left( \left( \frac{\mu}{1 - \mu} \right)^y (1 - \mu) \right) \right] = \exp \left[ y \log \left( \frac{\mu}{1 - \mu} \right) + \log(1 - \mu) \right] = \exp \left[ \frac{y \log \left( \frac{\mu}{1 - \mu} \right) + \log(1 - \mu)}{1} + 0 \right] \end{aligned}$$

- Видно, что  $\theta = \log \left( \frac{\mu}{1 - \mu} \right)$ ,  $a(\theta) = -\log(1 - \mu) \Rightarrow \mu = \frac{1}{1 + e^{-\theta}}$ ,  $a(\theta) = \log(1 + e^{\theta})$ .
- Тогда:  $a'(\theta) = \frac{e^{\theta}}{1 + e^{\theta}} = \frac{1}{1 + e^{-\theta}} = \sigma(\theta)$ ,  $\mathbb{E}(y | \mathbf{w}, \mathbf{x}) = g^{-1}(\langle \mathbf{w}, \mathbf{x} \rangle) = a'(\langle \mathbf{w}, \mathbf{x} \rangle) = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle)$ .

- Пусть  $\mu = \mathbb{E}[\langle \mathbf{x}, \mathbf{w} \rangle]$ . Рассмотрим байесовскую пуассоновскую регрессию:

$$p(y | \mu) = \frac{e^{-\mu} \mu^y}{y!}, \quad y \in \mathbb{N}.$$

- Здесь используется семейство распределений Пуассона, которое тоже принадлежит к экспоненциальному классу.

$$p(y | \mu) = \frac{e^{-\mu} \mu^y}{y!} = \exp[-\mu + y \log \mu - \log y!] = \exp \left[ \frac{y \log \mu - \mu}{1} - \log y! \right].$$

- Видно, что  $\theta = \log \mu$ ,  $a(\theta) = \mu \Rightarrow \mu = e^\theta$ ,  $a(\theta) = e^\theta$ .
- Тогда:  $a'(\theta) = e^\theta$ ,  $\mathbb{E}(y | \mathbf{w}, \mathbf{x}) = g^{-1}(\langle \mathbf{w}, \mathbf{x} \rangle) = a'(\langle \mathbf{w}, \mathbf{x} \rangle) = e^{\langle \mathbf{w}, \mathbf{x} \rangle}$ .

- Моделью смесей называют *вероятностную модель* для представления подгрупп объектов в общей группе без требований, чтобы информация о *принадлежности* к конкретной подгруппе явно содержалась в признаковом описании объекта.
- Модель соответствует смеси распределений, которое представляет собой распределение *вероятностей наблюдений* в общей совокупности.
- Часто модели смесей используются для описания мультимодальных данных (содержащих несколько регионов с высокой плотностью вероятности).
- Модели смесей используются для кластеризации (так называемая кластеризация на основе модели), а также для оценки плотности.



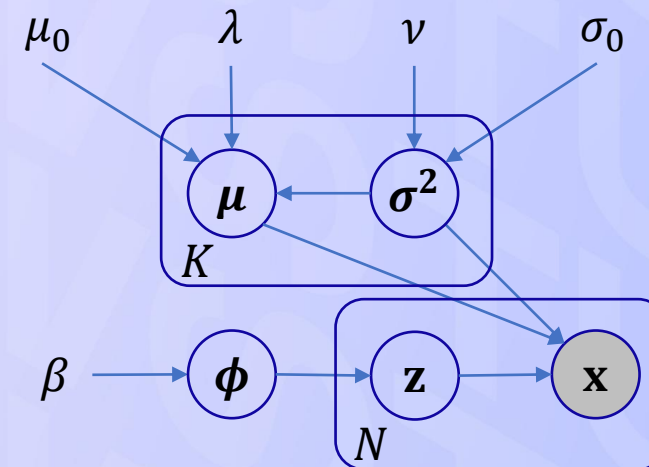
- Модель смеси обычно состоит из следующих *компонентов*:
  - $N$  наблюдаемых случайных переменных, каждая из которых имеет распределение смеси  $K$  компонент (все компоненты принадлежат к одному параметрическому семейству, но различаются *параметрами*).
  - $N$  латентных случайных переменных, определяющих компоненту из смеси для каждого *наблюдения*. Каждая переменная распределена по  $K$ -мерному категориальному распределению.
  - Множество  $K$  весов смеси, представляющих собой *вероятности*, сумма которых равна единице.
  - Множество  $K$  параметров, каждый из которых определяет параметр соответствующего компонента смеси.

- $\alpha$  – совместный гиперпараметр для параметров компонент смеси.
- $\beta$  – совместный гиперпараметр для весов смеси.
- $F(\mathbf{x} \mid \boldsymbol{\theta})$  – распределение вероятностей наблюдаемых переменных.
- $H(\boldsymbol{\theta} \mid \alpha)$  – априорное распределение параметров компонент смеси.
- $\theta_{i=1..K} \sim H(\boldsymbol{\theta} \mid \alpha)$  – параметр распределения наблюдаемой переменной, связанной с  $i$ -той компонентой смеси.
- $\phi_{i=1..K}$  – вес  $i$ -той компоненты смеси,  $\sum_i \phi_i = 1$ .
- $z_{i=1..N}$  – компонента, соответствующая  $i$ -той наблюдаемой переменной.
- $x_{i=1..N} \sim F(\boldsymbol{\theta}_{z_i})$  –  $i$ -тая наблюдаемая переменная.

Распределение  $F(\mathbf{x} \mid \boldsymbol{\theta})$  часто является нормальным или категориальным.

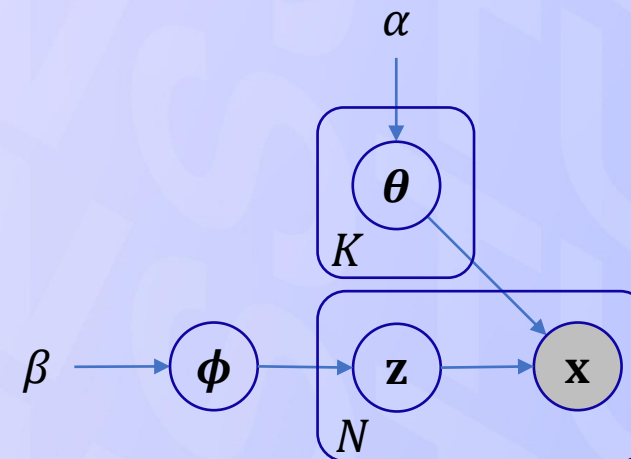
## Пример: смесь гауссовских моделей

- $\theta_{i=1..K} = \{\mu_{i=1..K}, \sigma_{i=1..K}^2\};$
- $\mu_{i=1..K} \sim \mathcal{N}(\mu_0, \lambda \sigma_i^2);$
- $\sigma_{i=1..K}^2 \sim \text{InvG}(\nu, \sigma_0^2);$
- $\Phi \sim \text{SymmDir}_K(\beta);$
- $z_{i=1..N} \sim \text{Cat}(\Phi);$
- $x_{i=1..N} \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i}^2).$



- *Вероятностный процесс устроен следующим образом. Сначала находится семпл весов компонент смеси  $\Phi$ . Затем семплируется компонента смеси  $z$ , и определяются параметры нормального распределения  $\mu$  и  $\sigma^2$ , которые соответствуют данной компоненте. Затем формируется значение  $x$ .*

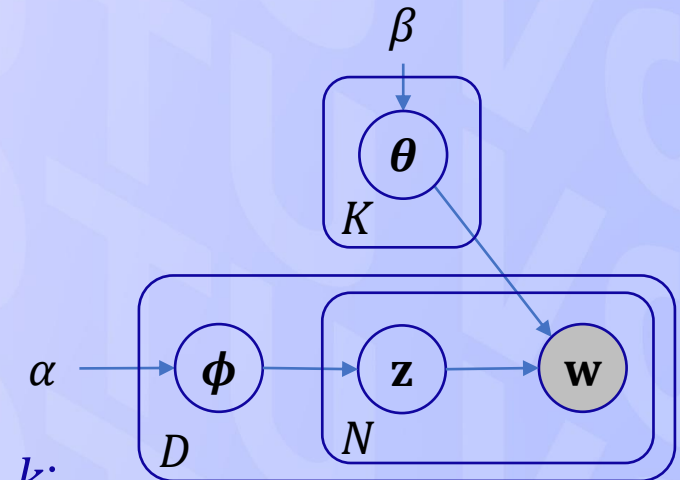
- $V$  – размерность категориальных наблюдаемых переменных;
- $\theta_{i=1..K} = \text{SymmDir}_V(\alpha)$ ;  $\theta_i \in \mathbb{R}^V$ ;  $\sum_j \theta_{ij} = 1$ ;
- $\phi \sim \text{SymmDir}_K(\beta)$ ;
- $z_{i=1..N} \sim \text{Cat}(\phi)$ ;
- $x_{i=1..N} \sim \text{Cat}(\theta_{z_i})$ ;
- *Вероятностный процесс устроен следующим образом. Сначала находится семпл весов компонент смеси  $\phi$ . Затем семплируется компонента смеси  $z$ , и определяются параметры категориального распределения  $\theta$ , которые соответствуют данной компоненте. Затем формируется значение  $x$ .*



- Модель латентного размещения Дирихле (LDA) относится к тематическим моделям, которые предназначены для описания текстов с точки зрения их тематик. LDA — это вероятностная модель *порождения* текста.
- Набор текстов (корпус) состоит из  $D$  документов. Каждый документ включает  $N$  слов.
- При этом *порядок употребления* терминов в документе игнорируется (модель «мешка слов»).
- Основное предположение тематической модели LDA состоит в том, что каждый документ имеет несколько тематик, смешанных в некоторой пропорции. Всего имеется  $K$  тематик. Тематика рассматривается как распределение вероятностей в пространстве слов из общего словаря.



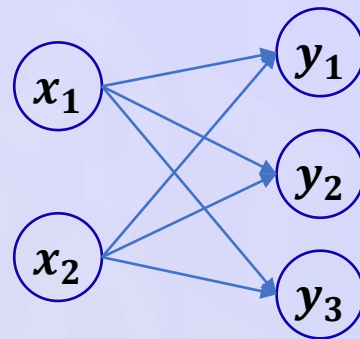
- $w \in \{1, 2, \dots, W\}$  – номер слова в словаре;
- $k \in \{1, 2, \dots, K\}$  – номер тематики;
- $d \in \{1, 2, \dots, D\}$  – номер документа в корпусе;
- $\mathbf{w}_d = [w_{d,1}, w_{d,2}, \dots, w_{d,N}]$  – слова в документе  $d$ ;
- $\mathbf{z}_d = [z_{d,1}, z_{d,2}, \dots, z_{d,N}]$  – тематики слов в документе  $d$ ;
- $\boldsymbol{\theta}_k = [\theta_{k,1}, \theta_{k,2}, \dots, \theta_{k,W}]$  – вероятности слов в тематике  $k$ ;
- $\boldsymbol{\phi}_d = [\phi_{d,1}, \phi_{d,2}, \dots, \phi_{d,K}]$  – вероятности тематик в документе  $d$ ;
- $\boldsymbol{\Phi} = [\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_D]^K \in \mathbb{R}^{D \times K}$  – вероятности тематик во всех документах;
- $\boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K]^W \in \mathbb{R}^{K \times W}$  – вероятности слов во всех тематиках.



- Возможные распределения компонент смеси:
  - Биномиальное распределение количества «положительных событий» (например, успехов, голосов «да» и т. д.) при фиксированном общем количестве вхождений.
  - Распределение Пуассона для количества повторений события за заданный период времени (для события, которое характеризуется фиксированной частотой возникновения).
  - Экспоненциальное распределение характеризует время до появления следующего события (для события, которое характеризуется фиксированной частотой возникновения).

- Возможные распределения компонент смеси:
  - Логарифмически нормальное распределение для положительных действительных чисел, которые, как предполагается, растут экспоненциально (например, доходы или цены).
  - Многомерное нормальное распределение (также известное как многомерное распределение Гаусса) для векторов коррелирующих событий, которые индивидуально распределены по Гауссу.
  - Многомерное  $t$ -распределение Стьюдента для векторов коррелированных результатов с тяжелым хвостом, и т.д.

- Искусственная нейронная сеть, в силу универсальной теоремы аппроксимации, позволяет воспроизвести произвольную функцию  $y = \Phi(x)$ .



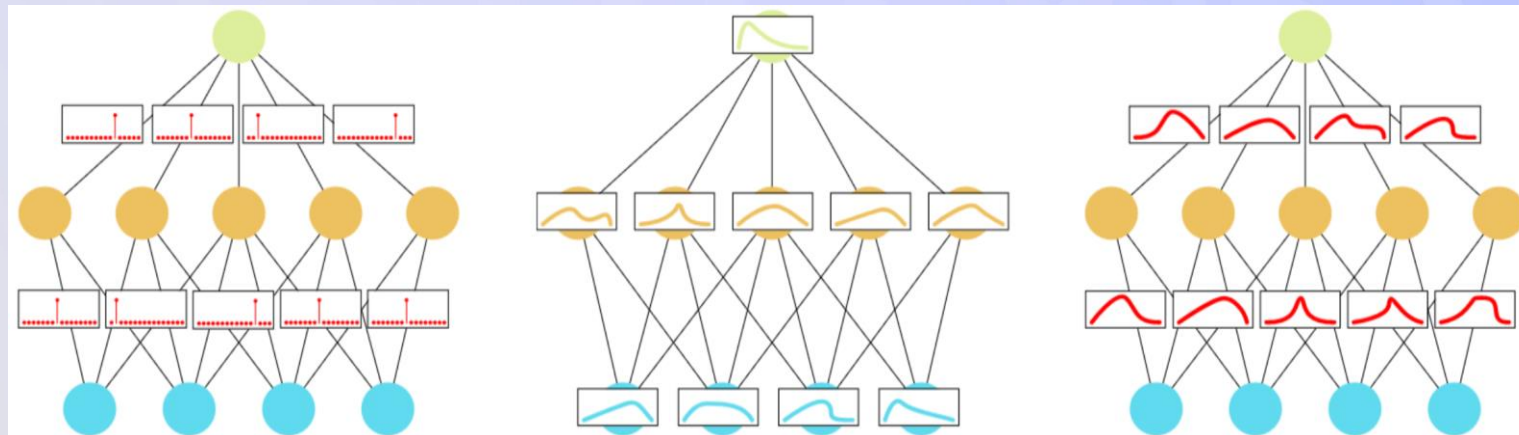
$$y_1 = \Phi(w_{11}x_1 + w_{12}x_2 + w_{10})$$

$$y_2 = \Phi(w_{21}x_1 + w_{22}x_2 + w_{20})$$

$$y = \Phi(Wx)$$

- В простейшей архитектуре сетей прямого распространения каждый слой представлен как линейное преобразование, за которым следует нелинейная операция  $\Phi$ , также известная как функция активации.
- Байесовские нейронные сети (BNN) представляют собой стохастические нейронные сети, обученные с использованием байесовского подхода.

- Стохастические нейронные сети — это тип нейросетей, построенный путем введения в сеть стохастических компонентов. Он реализуется путём введения либо стохастических активаций, либо стохастических весов для моделирования *нескольких* возможных моделей с соответствующим распределением вероятностей. Таким образом, BNN можно рассматривать как частный случай ансамблевого обучения.





- Первым шагом при разработке байесовской нейронной сети является выбор *архитектуры* глубокой сети, то есть функциональной модели:

$$\mathbf{y} = \Phi_{\mathbf{w}}(\mathbf{x}).$$

- Затем необходимо выбрать стохастическую модель, то есть *априорное распределение* возможных параметров модели и *априорную уверенность* в предсказательной силе модели:

$$p(\mathbf{w}), \quad p(\mathbf{y} | \mathbf{x}, \mathbf{w}).$$

- Выбор стохастической модели байесовской нейронной сети некоторым образом эквивалентен выбору функции потерь при обучении нейронной сети с точечной оценкой.

- В байесовском моделировании выделяют два основных вида неопределенности:
  - Алеаторическая неопределенность отражает шум, присущий наблюдениям. Это может быть, например, шум датчиков или шум движения. Данную неопределенность невозможно уменьшить, даже если собрать больше данных.
  - Эпистемическая неопределенность объясняет неопределенность в параметрах модели, которая отражает наше незнание того, как модель сгенерировала собранные данные. Эту неопределенность можно объяснить, имея достаточно данных, и ее часто называют неопределенностью модели.

- Для работы с байесовской нейросетью в Pyro нужно воспользоваться классом PyroModule:

```
- class Linear(nn.Module):  
+ class Linear(PyroModule):  
    def __init__(self, in, out):  
        super().__init__()  
        self.weight = ...  
        self.bias = ...  
    ...
```

```
- linear = Linear(5, 2)  
+ linear = PyroModule[Linear](5, 2)
```

```
linear = Linear(5, 2)  
assert isinstance(linear, nn.Module)  
assert not isinstance(linear,  
    PyroModule)  
  
to_pyro_module_(linear)
```

- Далее нужно заменить атрибуты *nn.Parameter* и *PyroParam* атрибутами *PyroSample*, которые задают априорные распределения:

```
linear.weight = PyroSample(dist.Normal(0, 1).expand([5, 2]).to_event(2))
```

- Непараметрические вероятностные модели – это класс статистических моделей, которые, в отличие от параметрических моделей, не предполагают фиксированной параметрической формы распределения данных.
- Такие модели обычно обладают большей гибкостью, поскольку они не ограничиваются заранее заданным набором параметров.
- **Случайный процесс** – это семейство случайных величин  $(\{X_t, t \in T\})$ , определенных на некотором вероятностном пространстве  $(\Omega, \mathcal{F}, P)$ .
- Случайный процесс можно представить как случайную функцию, которая для каждого фиксированного  $t$  выражает случайную величину, а для каждого фиксированного исхода  $\omega \in \Omega$  реализацию процесса (траекторию), т.е. функцию  $t \mapsto X(t, \omega)$ .

- Гауссов процесс является обобщением многомерного нормального распределения на бесконечное количество переменных.
- Он определяется как распределение над функциями  $f(\mathbf{x})$ , в котором для любого конечного набора точек  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  значения  $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$  имеют совместное многомерное нормальное распределение:

$$f \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \Leftrightarrow p(f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

- Гауссов процесс полностью определяется двумя характеристиками:
  - ✓ Функцией среднего  $m(\mathbf{x})$ , которая задает ожидаемое значение функции в точке  $\mathbf{x}$ :

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})].$$

- ✓ Ковариационной функцией (или ядром)  $k(\mathbf{x}, \mathbf{x}')$ , которая задает ковариацию между значениями функции в точках  $\mathbf{x}$  и  $\mathbf{x}'$ :

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))].$$



- Постановка **задачи регрессии**:

Дан обучающий датасет:  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ ,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ;  
 $\mathbf{y} = [y_1, y_2, \dots, y_n]$ ,  $y_i \in \mathbb{R}$ .

Требуется для данных из тестового датасета:  $\mathcal{D}_* = \mathbf{X}_* = [\mathbf{x}_{*1}, \mathbf{x}_{*2}, \dots, \mathbf{x}_{*m}]$ :

- *предсказать* значение  $\mathbf{y}_* = f(\mathbf{X}_*)$ ;
  - *оценить неопределенность* этого предсказания.
- Будем решать задачу, полагая, что связь между признаками и целевой переменной задаётся с помощью гауссова процесса:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')).$$

- Поскольку признаки всегда можно центрировать, то можно считать, что  $m(\mathbf{x}) = 0$ . Тогда гауссов процесс определяется только ядром  $k(\mathbf{x}, \mathbf{x}')$ .

- Ядро  $k(\mathbf{x}, \mathbf{x}')$  определяет свойства моделируемой функции  $f(\mathbf{x})$ , такие как гладкость, масштаб изменений и т.д.
- Одним из самых популярных ядер является радиально-базисное ядро (RBF, или квадратичное экспоненциальное ядро):

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \cdot \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right),$$

где  $\sigma_f^2$  – дисперсия (определяет амплитуду функции),  $l$  – параметр масштаба (определяет скорость изменения функции).

- Линейное ядро  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$  применяется для обычной линейной регрессии. Чтобы иметь возможность моделировать нелинейные зависимости, в гауссовских процессах используются нелинейные ядра, такие как RBF.

- Согласно определению гауссова процесса, совместное распределение значений функции в обучающих и тестовых точках имеет вид:

$$\begin{bmatrix} \mathbf{y}(\mathbf{X}) \\ \mathbf{y}_*(\mathbf{X}_*) \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right),$$

- ✓  $\mathbf{K}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{n \times n}$  – матрица ковариаций между всеми парами обучающих точек с элементами:  $[\mathbf{K}(\mathbf{X}, \mathbf{X})]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ ;
- ✓  $\mathbf{K}(\mathbf{X}, \mathbf{X}_*) \in \mathbb{R}^{n \times m}$  – матрица ковариаций между обучающими и тестовыми точками с элементами:  $[\mathbf{K}(\mathbf{X}, \mathbf{X}_*)]_{ij} = k(\mathbf{x}_i, \mathbf{x}_{*j})$ ;
- ✓  $\mathbf{K}(\mathbf{X}_*, \mathbf{X}) = \mathbf{K}(\mathbf{X}, \mathbf{X}_*)$ ,  $\mathbf{K}(\mathbf{X}_*, \mathbf{X}) \in \mathbb{R}^{m \times n}$ ;
- ✓  $\mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \in \mathbb{R}^{m \times m}$  – матрица ковариаций между всеми парами тестовых точек с элементами:  $[\mathbf{K}(\mathbf{X}_*, \mathbf{X}_*)]_{ij} = k(\mathbf{x}_{*i}, \mathbf{x}_{*j})$ ;
- ✓  $\sigma^2 \mathbf{I}$  – матрица шума, добавляемая к ковариационной матрице.

- Для предсказания значений  $\mathbf{y}_*$  в тестовых точках нужно найти условное распределение:  $p(\mathbf{y}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_*)$ .
- Из свойств многомерного нормального распределения известно, что условное распределение также является нормальным.
- Можно показать, что для полученного совместного распределения  $p(\mathbf{y}_*, \mathbf{y})$  условное распределение  $p(\mathbf{y}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_*)$  выражается следующим образом:

$$p(\mathbf{y}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_*) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*),$$

а среднее значение и ковариационная матрица задаются выражениями:

$$\boldsymbol{\mu}_* = \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}]^{-1} \mathbf{y};$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*);$$

- Это выражение и задаёт распределение предсказанных значений  $\mathbf{y}_*$ .

- Полученные значения  $\mu_*$  и  $\Sigma_*$  можно интерпретировать таким образом:
  - ✓ Среднее значение  $\mu_*$  является наилучшим предсказанием  $y_*$ . Если требуется точечное предсказание, то достаточно получить  $\mu_*$ .
  - ✓ Ковариационная матрица содержит информацию о:
    - Дисперсии – диагональные элементы матрицы ковариации соответствуют дисперсии предсказания в каждой тестовой точке. Среднеквадратичное отклонение предсказания в точке  $\mathbf{x}_{*i}$ :
$$\sigma(\mathbf{x}_{*i}) = \sqrt{[\Sigma_*]_{ii}},$$
Это позволяет строить доверительные интервалы вокруг  $\mu_*$ .
    - Ковариации – недиагональные элементы  $\Sigma_*$  показывают, насколько предсказания в разных тестовых точках *коррелируют* друг с другом.



- Приведённая функция содержит неизвестные: гиперпараметры ядра ( $\sigma_f^2$  и  $l$ ) и дисперсию шума ( $\sigma^2$ ). Их требуется определить из данных.
- Для этого используется метод максимального правдоподобия.
- Так как  $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{I})$ , то логарифм правдоподобия:

$$\mathcal{L} = \log p(\mathbf{y} | \mathbf{X}, \sigma_f^2, l, \sigma^2) = -\frac{1}{2} \mathbf{y}^T [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{I}]^{-1} \mathbf{y} - \frac{1}{2} \log \det(\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{I}) - \frac{n}{2} \log(2\pi).$$

- Максимизация логарифма правдоподобия эквивалентна минимизации отрицательного логарифма правдоподобия по гиперпараметрам:

$$\sigma_f^{2*}, l^*, \sigma^{2*} = \arg \min_{(\sigma_f^2, l, \sigma^2) \in \Theta} [-\log p(\mathbf{y} | \mathbf{X}, \sigma_f^2, l, \sigma^2)].$$

- Решение данной задачи оптимизации позволит получить значения гиперпараметров, которые используются в выражениях для  $\boldsymbol{\mu}_*$  и  $\boldsymbol{\Sigma}_*$ .

- Хотя гауссовы процессы являются весьма мощными моделями, у них есть несколько важных ограничений:
  1. Вычислительная сложность. Обучение модели требует инверсии матрицы  $\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}$ , алгоритм которой имеет сложность  $\mathcal{O}(n^3)$ . Получение предсказаний для  $m$  тестовых точек имеет сложность  $\mathcal{O}(nm + m^2)$ , что ограничивает использование гауссовых процессов для больших наборов данных.
  2. Чувствительность к выбору ядра. Результаты сильно зависят от выбора ядра и его гиперпараметров. Неправильный выбор может привести к плохой аппроксимации данных.

3. Предположение о нормальности. Гауссовы процессы предполагают, что данные подчиняются нормальному распределению. Если шум в данных не является нормальным (например, имеет тяжелые хвосты), модель может давать некорректные результаты.
  4. Сложности с экстраполяцией. Гауссовы процессы плохо работают при экстраполяции за пределы области обучающих данных, особенно с RBF-ядром, так как предсказания стремятся к среднему значению (обычно нулю) с увеличением расстояния от обучающих точек.
-

- Пусть дана параметрическая модель (например, *нейронная сеть, регрессия или градиентный бустинг*) с гиперпараметрами  $\lambda \in \Lambda$ .
- Здесь  $\Lambda$  – пространство гиперпараметров (например, скорость обучения, количество слоев, регуляризационные коэффициенты и т.д.).
- Требуется найти такие гиперпараметры  $\lambda^*$ , которые минимизируют некоторую целевую функцию  $f(\lambda)$ , например, ошибку на валидационной выборке:

$$\lambda^* = \arg \min_{\lambda \in \Lambda} f(\lambda).$$

- При этом предполагается, что поиск параметров  $w$  параметрической модели  $f(\lambda) = f_w(\lambda)$  осуществляется в процессе обучения модели для фиксированных значений  $\lambda$ .

- При определении гиперпараметров следует иметь в виду:
  1. Функция  $f(\lambda)$  обычно не имеет аналитического вида (вычисление её значений возможно только после обучения модели  $f_w(\lambda)$  и оценки её производительности).
  2. Расчёт  $f(\lambda)$  может быть вычислительно сложным (например, обучение нейронной сети).
  3. Пространство гиперпараметров  $\Lambda$  может быть сложным (непрерывным, дискретным или смешанным), а сама функция  $f(\lambda)$  – мультимодальной, шумной и недифференцируемой.
- Для решения такой задачи используется **байесовская оптимизация**, эффективно выполняющая глобальную оптимизацию «черного ящика».



- Пусть целевая функция  $f(\lambda) \sim \mathcal{GP}(m(\lambda), k(\lambda, \lambda'))$ .
- Здесь  $m(\lambda)$  – функция среднего значения (обычно принимается равной нулю), а  $k(\lambda, \lambda')$  – ковариационная функция (например, ядро *RBF*).
- На каждом шаге  $t$  байесовской оптимизации имеется набор уже вычисленных значений функции:

$$\mathcal{D}_t = \{(\lambda_i, f(\lambda_i))\}_{i=1}^t$$

- Используя GP, можно вычислить апостериорное распределение  $f(\lambda)$  в любой точке  $\lambda$ :

$$p(f(\lambda) | \mathcal{D}_t) = \mathcal{N}(\mu_t(\lambda), \sigma_t^2(\lambda)),$$

где среднее и дисперсия вычисляются аналогично формулам регрессии.

$$\mu_t(\lambda) = \mathbf{k}_t(\lambda)^T [\mathbf{K}_t + \sigma_n^2 \mathbb{I}]^{-1} \mathbf{f}_t;$$
$$\sigma_t^2(\lambda) = k(\lambda, \lambda) - \mathbf{k}_t(\lambda)^T [\mathbf{K}_t + \sigma_n^2 \mathbb{I}]^{-1} \mathbf{k}_t(\lambda).$$

- Здесь:
  - $\mathbf{f}_t = [f(\lambda_1), f(\lambda_2), \dots, f(\lambda_t)]^T$  – вектор значений целевой функции;
  - $[\mathbf{K}_t]_{ij} = k(\lambda_i, \lambda_j)$  – ковариационная матрица размером  $t \times t$ ;
  - $\mathbf{k}_t(\lambda) = [k(\lambda, \lambda_1), k(\lambda, \lambda_2), \dots, k(\lambda, \lambda_t)]^T$  – вектор ковариаций между точкой  $\lambda$  и точками из  $\mathcal{D}_t$ ;
  - $k(\lambda, \lambda)$  – априорная дисперсия в точке  $\lambda$ ;
  - $\sigma_n^2$  – дисперсия шума.
- Эти формулы позволяют не только предсказать значение  $f(\lambda)$  в новой точке  $\lambda$ , но и оценить неопределенность предсказания через  $\sigma_t^2(\lambda)$ .

- На каждой итерации байесовской оптимизации требуется выбрать следующую точку  $\lambda_{t+1}$ , в которой будет вычислена целевая функция  $f(\lambda)$ .
- Но вычисление целевой функции вычислительно сложно, поэтому перебор всех возможных точек не является эффективной стратегией.
- Введём функцию принятия решений  $a(\lambda)$ , которая будет оценивать некоторую «полезность» вычисления  $f(\lambda)$  в данной точке.
- Тогда оптимальная следующая точка определяется выражением:

$$\lambda_{t+1} = \arg \max_{\lambda \in \Lambda} a(\lambda).$$

- Функция  $a(\lambda)$  зависит от апостериорного распределения  $p(f(\lambda) | \mathcal{D}_t)$ , т.е., от величин  $\mu_t(\lambda)$  и  $\sigma_t^2(\lambda)$ . Существует несколько популярных вариантов функций принятия решений.

- Одна из функций принятия решений – ожидаемое улучшение (Expected Improvement, EI).
- Пусть  $f_{best}$  – текущее лучшее значения целевой функции:

$$f_{best} = \min_{i=1,2,\dots,t} f(\lambda_i).$$

- Улучшение в точке  $\lambda$  определяется выражением:

$$I(\lambda) = \max(f_{best} - f(\lambda), 0).$$

- Так как  $f(\lambda)$  – случайная величина, то вычисляем ожидаемое улучшение:

$$a_{EI}(\lambda) = \mathbb{E}[I(\lambda)] = \mathbb{E}[\max(f_{best} - f(\lambda), 0)].$$

- Для нормального распределения ожидаемое улучшение может быть выражено аналитически.

- Другая функция принятия решений – вероятность улучшения (*Probability of Improvement, PI*).
- Она измеряет *вероятность* того, что  $f(\boldsymbol{\lambda})$  будет лучше, чем  $f_{best}$ :

$$a_{PI}(\boldsymbol{\lambda}) = P(f(\boldsymbol{\lambda}) < f_{best}) = \Phi\left(\frac{f_{best} - \mu_t(\boldsymbol{\lambda})}{\sigma_t(\boldsymbol{\lambda})}\right).$$

- Здесь  $\Phi(\cdot)$  – функция распределения стандартного нормального распределения.
- Величина PI проще в вычислении, но менее эффективна, чем EI, так как она не учитывает величину улучшения, а только вероятность его достижения.



- После выбора функции принятия решений нужно найти:

$$\lambda_{t+1} = \arg \max_{\lambda \in \Lambda} a(\lambda).$$

- Данная задача оптимизации ищет максимум функции принятия решений, которая, в отличие от целевой функции, относительно легко вычисляется, поскольку зависит только от величин  $\mu_t(\lambda)$  и  $\sigma_t^2(\lambda)$ . А они вычисляются аналитически через GP.
- Для оптимизации обычно используются градиентный спуск, L-BFGS и т.д.
- После получения  $\lambda_{t+1}$  нужно вычислить  $f(\lambda_{t+1})$  и добавить результат в  $\mathcal{D}_{t+1}$ , а затем повторить процесс до достижения критерия остановки.
- В результате получаются оптимальные гиперпараметры  $\lambda^*$ :

$$\lambda^* = \arg \min_{\lambda_i \in \{\lambda_1, \lambda_2, \dots, \lambda_T\}} f(\lambda_i).$$

# Демонстрация практических примеров

---

## Заключение

1. Поговорили про недостатки линейных моделей, узнали про обобщенные линейные модели, ввели понятие функции связи.
2. Вспомнили экспоненциальное семейство распределений, выяснили, как составлять обобщенную линейную модель для различных распределений.
3. Разобрали примеры обобщенных линейных моделей.
4. Ввели понятие моделей смесей, определили структуру моделей смесей.
5. Рассмотрели часто используемые примеры моделей смесей.
6. Ввели понятия байесовской и стохастической нейронных сетей.
7. Рассмотрели непараметрические модели, поговорили о гауссовых процессах и основных задачах, решаемых с их помощью.
8. На практических примерах посмотрели реализацию моделей в Pyro.

# **Спасибо за внимание!**

Волгоград 2025

---