

# 卒研ゼミ 「深層学習」

## 6.2 Gradient-Based Learning

### 6.2.1 Cost Function

#### **6.2.1.1 Learning Conditional Distribution with Maximum Likelihood**

#### **6.2.1.2 Learning Conditional Statistics**

### 6.2.2 Output Units

#### 6.2.2.1 Linear Units for Gaussian Output Distributions

#### 6.2.2.2 Sigmoid Units for Bernoulli Output Distributions

#### 6.2.2.3 Softmax Units for Multinoulli Output Distributions

#### 6.2.2.4 Other Output Types

**2019年12月9日？**

**岡崎健人**

## 6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

現代におけるNNのほとんどは最尤推定を用いている

$$J(\theta) = - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{y} | \mathbf{x})]$$

- 最尤推定では、コスト関数は「負の対数尤度」になる
- 訓練データとモデル分布の交差エントロピーと考えることもできる

## ※復習：最尤推定 (§5.5)

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})]$$

$p_{\text{data}}(\mathbf{x})$  : 真であるが  $\boldsymbol{\theta}$  がわからないため未知の分布

$p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  : 値が  $\mathbf{x}$  であるデータを得る確率 ( $\boldsymbol{\theta}$  が変数の関数)

$\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  : i. i. d. に従うデータ

$\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\cdot]$  : 確率変数  $\mathbf{x}$  が経験的な分布  $\hat{p}_{\text{data}}$  に従うときの期待値

## ※復習：交差エントロピー (§3.13)

$$\begin{aligned} H(\hat{p}_{\text{data}}, p_{\text{model}}) &= H(\hat{p}_{\text{data}}) + D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) \\ &= -\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})] \end{aligned}$$

**注意！** ここでは訓練データをまとめて1つの文字  $\mathbf{x}$  で表しているが、教科書の該当範囲では  $\mathbf{x}$  と  $\mathbf{y}$  が訓練データ

## 6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

計算例：  $J(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \boldsymbol{\theta})\|^2 + \text{const.}$

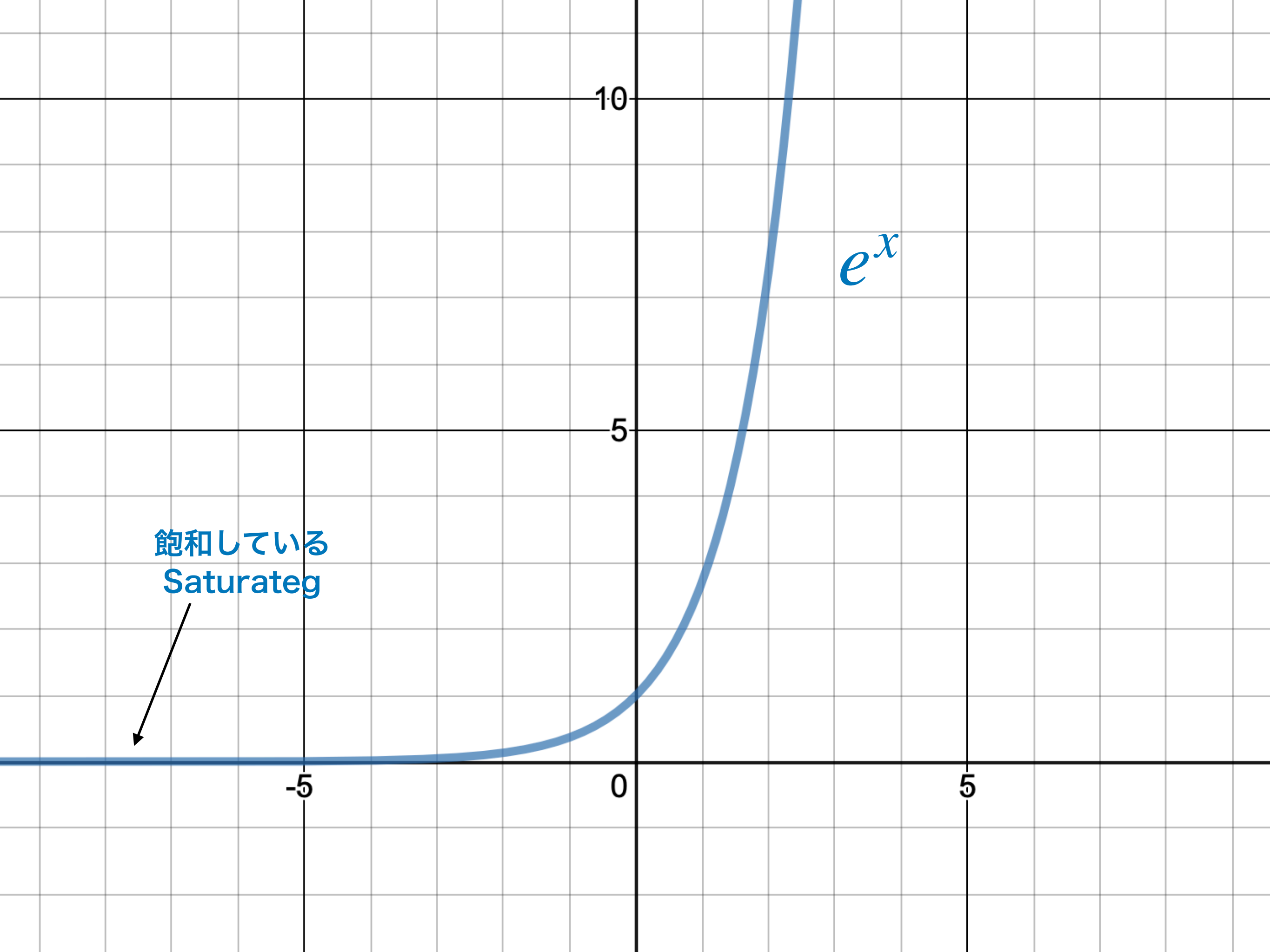
- ・ 尤度関数が  $\mathcal{N}(\mathbf{y}; f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{I})$  に従う場合
- ・ 線形モデルのとき、最尤推定と最小二乗法は等価 (§5.5.1)

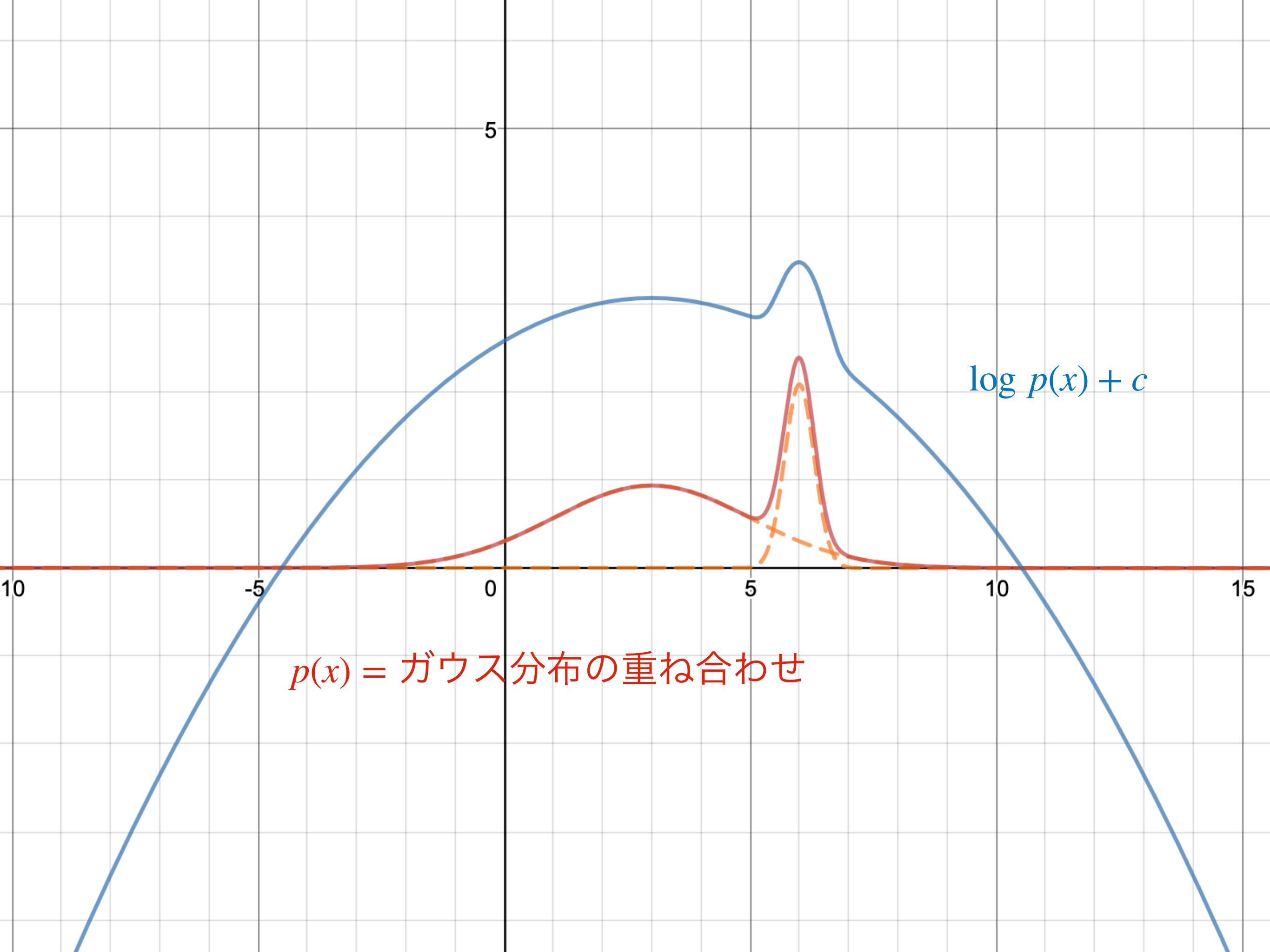
コスト関数を自動的に作り出せるのが、最尤推定の強み

## 6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

コスト関数の勾配  $|\nabla J(\theta)|$  は大きく予測可能であってほしい

- 勾配降下法では  $\nabla J(\theta)$  を用いて最小点を求めるため
- **飽和する** (saturate, become very flat) ような関数は使えない
- 活性化関数が飽和するとよくない
- 尤度関数が  $e^x$  を含んでいても、負の対数尤度は「飽和」しない







## 6.2.1.1 Learning Conditional Distributions with Maximum Likelihood

- 実際によく用いられるモデル(?)に対して、 $J(\theta) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{y} | \mathbf{x})]$  は、通常、最小値をもたない
- 連続型確率変数の出力分布の密度を極めて高くできる
- →交差エントロピーは $-\infty$ に近づく
- 無限大の報酬を回避するように、正則化を行う (§7)

## 6.2.1.2 Learning Conditional Statistics

分布 $p(y | x; \theta)$ よりも「 $x$ が与えられた元での $y$ の推定量」を知りたいとき

- $y$ の平均を与える関数  $f(x; \theta)$  は知っているものとする
- 十分強力なNNを用いれば、どんな関数  $f$  も表現できる
- この考えは**汎関数**(functional, 関数から実数値へのマッピング)になる
- コスト関数が「ある関数形」をとれば、交差エントロピーが最小値を持つようになる (かも)
- **変分法**(calculus of variations)が必要だが**次のことを知っとけばOK**

## 6.2.1.2 Learning Conditional Statistics

知っとくこと①/2

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|^2$$

$$\rightarrow f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y}|\mathbf{x})}[\mathbf{y}]$$

無限に多くのデータを真の分布から得るならば、コスト関数は平均二乗誤差になり、 $\mathbf{x}$  の値が与えられれば対する  $\mathbf{y}$  平均を知ることができる

## 6.2.1.2 Learning Conditional Statistics

知っとくこと②/2

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \|\mathbf{y} - f(\mathbf{x})\|_1$$

→  $f^*(\mathbf{x}) = (\mathbf{x} \text{ に対する } \mathbf{y} \text{ の中央値})$

- $\mathbf{x}$  の値が与えられたときの  $\mathbf{y}$  の**中央値**(median)が得られる
- このコスト関数は**平均絶対誤差**(mean absolute error)とよばれる

## 6.2.1.2 Learning Conditional Statistics

それでも $p(y | x)$ を求めたほうがいい

- これらのコスト関数を勾配法で使うと、出力ユニットが飽和してうまくいかないことがある
- そのため条件付き分布 $p(y | x)$ を求める必要がなくても、交差エントロピーの方法を使うほうがポピュラー