# Lecture 2

- Part 1  Regression

  - Supervised learning

| | 1 | 2 | ...$j$... | $p$ | 1 | 2 | ...$j$... | $d$ |
|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | |
| 2 | | | | | | | | |
| ...$i$... | | $X_:^T$ | | | | $y_:^T$ | | |
| $m$ | | | | | | | | |

  $\underbrace{\qquad\qquad\qquad\qquad}_{\text{Training Data}}$

- Notation:
  - Specific example:
    $(x_:, y_:)$

$$x_i = \begin{pmatrix} x_{:,1} \\ \vdots \\ x_{:,j} \\ \vdots \\ x_{:,p} \end{pmatrix} \qquad y_i = \begin{pmatrix} y_{:,1} \\ \vdots \\ y_{:,j} \\ \vdots \\ y_{:,d} \end{pmatrix}$$
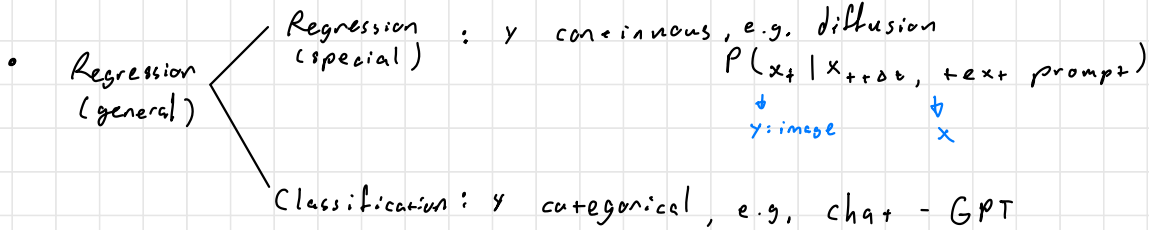
  - generic example

    $(x, y)$        drop subscript $i$

$$x \begin{pmatrix} x_1 \\ \vdots \\ x_j \\ \vdots \\ x_p \end{pmatrix} \qquad y = \begin{pmatrix} y_1 \\ \vdots \\ y_j \\ \vdots \\ y_d \end{pmatrix}$$

- Regression
  (general)

  Regression (special): $y$ continuous, e.g. diffusion
  $$P(x_t \mid x_{t+\Delta t}, \text{text prompt})$$
  $\underset{y:\text{image}}{\downarrow} \qquad \underset{x}{\downarrow}$

  Classification: $y$ categorical, e.g. chat - GPT
  $$P(x_t \mid x_{<t})$$
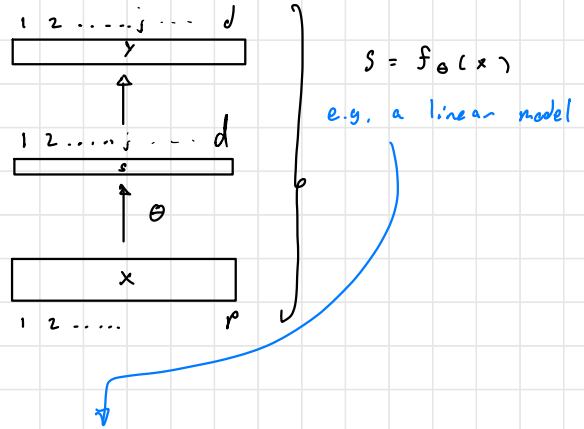  $\underset{y}{\downarrow} \qquad \underset{x}{\downarrow}$
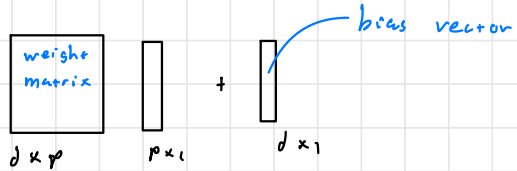  $\in \{\text{50k categories (tokens)}\}$

- Gauss Paradigm :
  1. probabilistic formulation
  2. objective function / loss
  3. learning algorithm
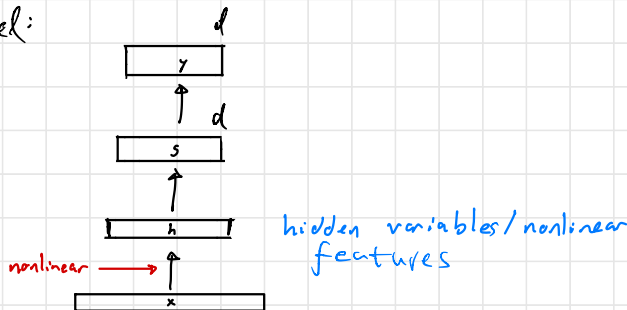  4. experiments, theory

- Unification :

  Start with $P(y|x)$

$$S = f_\theta(x)$$

e.g. a linear model

$$f_\theta(x) = wx + b$$

weight matrix $d \times p$    $p \times 1$    $+$    bias vector $d \times 1$

$$\theta = (w, b)$$

Another model:

nonlinear

hidden variables/nonlinear features

$$f_\theta(x) = wh + b$$

general formulation

, Regression

1. $P_\theta(y \mid x) \sim \mathcal{N}(s, \sigma^2 I)$

   ie. $y_j \sim \mathcal{N}(s_j, \sigma)$ , $j = 1, \ldots, d$

   $$P_\theta(y \mid x) = \prod_{j=1}^{d} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_j - s_j)^2}{2\sigma^2}} \qquad \text{\textcolor{red}{by independence}}$$

   $$= \frac{1}{(\sqrt{2\pi\sigma^2})^d} \exp\left(-\frac{1}{2\sigma^2} |y - s|^2\right).$$

2. Log - Likelihood :

   $$\ell(\theta) = \log P_\theta(y \mid x) \qquad \text{\textcolor{red}{objective for single example}}$$

   $$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \log P_\theta(y_i \mid x_i)$$

   $$= \mathbb{E}_{\text{data}}(\log P_\theta(y \mid x))$$

   - data distribution: $P_{\text{data}}(x, y)$

   data = 
   $$(x_i, y_i) \sim P_{\text{data}}(x, y) , \quad i = 1, \ldots, n, \quad \text{indep}$$

   - empirical data distribution

   $$\hat{P}_{\text{data}} = \frac{1}{n} \sum_{i=1}^{n} d_{x_i, y_i} \sim \text{\textcolor{blue}{Unif}} \left\{(x_i, y), i = 1 \cdots n\right\}$$

- $\ell(\theta) = \log P_\theta(y \mid x)$

  $\quad = -\dfrac{1}{2\sigma^2} \, |y - s|^2 + \text{const}$

- $L(\theta) = -\dfrac{1}{2\sigma^2} \displaystyle\sum_{i=1}^{n} |y_i - s_i|^2 + \text{const}$

- Notice:

  $$\max_{\theta} L(\theta) = \min_{\theta} \sum_{i=1}^{n} |y_i - f_\theta(x_i)|^2$$

  $\underbrace{\qquad\qquad\qquad\qquad}_{}$

  Least - Squares Loss

- $\text{Loss}(\theta) = \dfrac{1}{2}\displaystyle\sum_{i=1}^{n} |y_i - s_i|^2$ , $s_i = f_\theta(x_i)$ , e.g $f_\theta(x_i) = W x_i + b$

- Generic Example:

  $$\text{loss}(\theta) = |y - s|^2 = \dfrac{1}{2}\sum_{j=1}^{d} (y_j - s_j)^2$$

  $\qquad\qquad\qquad$ drop i

3. $\min\limits_{\Theta}$ Loss $(\Theta)$ by gradient-based method

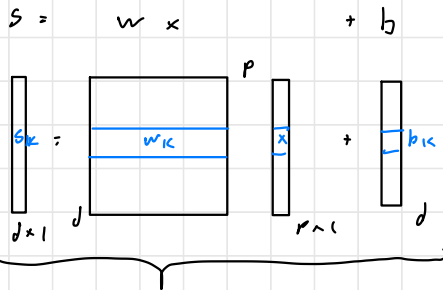$$\frac{\partial\, loss(\Theta)}{\partial s_k} = -(Y_k - S_k) \underset{(error)}{=} - e_k$$

$$\frac{\partial l}{\partial s} = \begin{pmatrix} \vdots \\ \frac{\partial l}{\partial s_k} \\ \vdots \end{pmatrix}_d^k = \begin{pmatrix} \vdots \\ -e_k \\ \vdots \end{pmatrix}_d^k = -\begin{pmatrix} \vdots \\ Y_k - S_k \\ \vdots \end{pmatrix}_d^k$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\quad e \qqu\qquad Y - S$$

$$\frac{\partial l}{\partial s} = -(Y - S) = -e$$

now we want the derivative with respet to $\Theta$.

$$S = \quad w\; x \qquad\qquad + b$$



$$s_k = w_k X + b_k$$

$$\frac{\partial l}{\partial w_k} = \frac{\partial l}{\partial s_k}\,\frac{\partial s_k}{\partial w_k} \longrightarrow \left(\frac{\partial l}{\partial w_k}\right)_d^k = -\left(e_k\right)_d^k X^T$$

$$= -e_k\, X^T$$



$$\frac{\partial l}{\partial b} = -e$$

- Classification:

individual categories
↓
$$y \in \{1, 2, \ldots, C = 50k\}$$

$$y = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \\ c \\ \\ \\ \end{matrix} \quad C = 50k$$

← one-hot vector

equivalent notation

1. $P(y \mid x)$

$$P(y = c \mid x) = P_c \qquad \sum_{i=1}^{C=50k} P_c = 1$$

$$P\left( y = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} 1 \\ \\ c \\ \\ \end{matrix} \; C=50k \;\middle|\; x \right) = P_c \qquad P = \begin{pmatrix} P_1 \\ \vdots \\ P_C \end{pmatrix}$$

One hot

| | y | | |
|1|2|...| c |

$C=50k$     probs

| P |
$C$

soft-max →     logits

| s |
|1|2|...|C|     $C$      $\longrightarrow$     $s = wx + b$

| x |
|1|2| ... |     $P$

$$\begin{bmatrix} \\ \\ \end{bmatrix} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix} \begin{bmatrix} \\ \\ \end{bmatrix} + \begin{bmatrix} \\ \\ \end{bmatrix}$$
$C \times 1$     $C \times P$     $P \times 1$     $C \times 1$

▷ Softmax:

or

$$P_{1c} = \frac{e^{s_k}}{\sum_{c=1}^{C=50k} e^{s_c}} = \frac{e^{s_k}}{Z}$$

$$s = wh + b$$
↑ non-linear
$x$

- Introduce temperature $T > 0$

$$p_k = \frac{e^{s_k/T}}{\sum_c e^{s_c/T}}$$

$$p = \begin{pmatrix} p_1 \\ \vdots \\ p_k \\ \vdots \\ p_c \end{pmatrix}$$

$T \to \infty$  Unif $\{1, \dots, C\}$

$T \to 0$  hardmax

$\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$ ← if $s_k = \max_c s_c$

one-hot

- note one degree of redundancy: (not an issue if $C \gg 2$)

$$p_k = \frac{e^{s_k}}{\sum_c e^{s_c}} = \frac{e^{s_k + \text{const}}}{\sum_c e^{s_c + \text{const}}}$$

2. If $y = k$, $p(y|x) = \dfrac{e^{s_k}}{Z} = \dfrac{e^{\langle y, s \rangle}}{Z}$

$\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ C \end{matrix}$ one-hot   $s = \begin{pmatrix} s_1 \\ \vdots \\ s_k \\ \vdots \\ s_c \end{pmatrix}$ logits

$\ell(\theta) = \log p(y|x)$

$\qquad = s_k - \log Z$

$\qquad = \langle y, s \rangle - \log Z$

3. $\text{loss}(\Theta) = -\log p_{\theta}(y|x)$

$$= -(s_k - \log z)$$

$$= -(\langle y, s \rangle - \log z)$$

$\dfrac{\partial \text{loss}(\Theta)}{\partial s_k} = -\left( y_k - \dfrac{\partial}{\partial s_k} \log z \right)$ $\qquad \langle y, s \rangle = \displaystyle\sum_{c=1}^{C} y_c s_c$

$$= -\left( y_k - \dfrac{1}{z}\dfrac{\partial z}{\partial s_k} \right), \qquad z = \displaystyle\sum_{c=1}^{C} e^{s_c}$$

$$= -\left( y_k - \dfrac{1}{z} e^{s_k} \right)$$

$$= -(\underbrace{y_k - p_k}_{e_k})$$

$\dfrac{\partial \text{loss}(\Theta)}{\partial s} = -(\underbrace{y - p}_{e})$

cross-entropy loss $= E_{\text{data}}\left[ -\log p_{\theta}(y|x) \right]$

cross-entropy of $q$ relative to $p$ $= H(p, q) = -E_p[\log q]$

- Case   C = 2

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \begin{matrix} + \\ - \end{matrix} \quad , \quad p = \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} \begin{matrix} + \\ - \end{matrix} \quad , \quad s = \begin{pmatrix} s_1 \\ s_2 \end{pmatrix} \begin{matrix} + \\ - \end{matrix}$$

$$p_1 = \frac{e^{s_1}}{e^{s_1} + e^{s_2}} \quad , \quad p_2 = \frac{e^{s_2}}{e^{s_1} + e^{s_2}}$$

- Now we fix $s_2 = 0$, so $\quad p_1 = \dfrac{e^{s_1}}{1 + e^{s_1}} = $ sigmoid $(s_1)$
  (due to redundancy)

$$p_2 = 1 - p_1 = \frac{1}{1 + e^{s_1}}$$

$$P(Y_1 = 1 \mid x) = p_1 = \frac{e^{s_1}}{1 + e^{s_1}} = \text{sigmoid}(s_1) \qquad \text{Special case of softmax}$$

$$P(y \mid x) = \frac{e^{\langle y, s \rangle}}{Z} = \frac{e^{y_1 s_1 + y_2 s_2}}{Z} = \frac{e^{y_1 s_1}}{Z}$$

$$\log P(y_1 \mid x) = y_1 s_1 - \log(1 + e^{s_1})$$

$$y_1 \in \{1, 0\}$$
$\uparrow$
enough to indicate the two categories