# Deep learning.

- Model parametrizati
- Defining a loss functi
- Optimizati
- Model selecti
- Reporting

---

### Parametrizati

$$p(y \mid x)$$

$$\underset{\text{output}}{\sim} \quad \underset{\text{input}}{-}$$

(1) $p(y \mid x) \propto \exp\{F(x,y)\}$

(2) latent variable approach

$$p(y,z \mid x) \propto \exp\{F(x,y,z)\}$$
$$p(z \mid x) \propto \exp\{F'(x,z)\}$$

---

### Multiclass classificati

$$y \in \{1, \cdots, c\}$$

$$F(x,y) = w_y^T \phi(x) + b_y$$

$$\uparrow \qquad \uparrow \qquad \uparrow$$
$$\mathbb{R}^d \qquad \mathbb{R}^d \qquad \mathbb{R}$$

$$p(y \mid x) = \frac{\exp\{w_y^T \phi(x) + b_y\}}{\sum_{y'=1}^{c} \exp\{w_{y'}^T \phi(x) + b_y\}}$$

## Mixture density network

$$y \in \mathbb{R}^{d'}$$
$$z \in \{1, \cdots, c\}$$
$$F(x, y, z) = -\| y - W_\mu^z \phi(x) \|^2 \quad \phi \in \mathbb{R}^{d'}, \quad W \in \mathbb{R}^{d' \times d}, \quad \in \mathbb{R}^d$$
$$F'(x, z) = W_z^T \phi(x) + b_z \quad \underset{\mathbb{R}^d}{\uparrow} \quad \underset{\mathbb{R}}{\uparrow}$$

$$p(y|x) = \sum_{z=1}^{c} p(y, z | x) = \sum_{z=1}^{c} p(y | z, x) \, p(z|x)$$

$$= \sum_{z=1}^{c} \left( \frac{\exp\{ W_z^T \phi(x) + b_z \}}{\sum_{z'=1}^{c} \exp\{ W_{z'}^T \phi(x) + b_{z'} \}} \right) \left( \frac{\exp\{ -\| y - W_\mu^z \phi(x) \|^2 \}}{\int_{\mathbb{R}^{d'}} \exp\{ -\| y' - W_\mu^z \phi(x) \|^2 \} dy' } \right)$$

---

## Loss function

for $(x, y^*)$, $\underline{\ell}$ is $-\log p(y^* | x)$ $\leftarrow$ per-example loss

for $D = \{ (x_1, y_1^*), \cdots, (x_N, y_N^*) \}$,

$$p(D) = \prod_{n=1}^{N} p(x_n, y_n^*) = \prod_{n=1}^{N} p(y_n^* | x_n) \underbrace{p(x_n)}_{\text{uniform}}$$

$$\log p(D) = \sum_{n=1}^{N} \log p(y_n^* | x_n) + \text{Const.}$$

---

## Multiway classification

$$-\log p(y^* | x) = \underbrace{-W_{y^*}^T \phi(x) - b_{y^*}}_{= F(x, y)} + \log \sum_{y'=1}^{c} \exp\{ W_{y'}^T \phi(x) + b_{y'} \}$$

## Optimization: Gradient descent

$L(D; \theta)$ : a data-level loss fnt.

$$L(D; \theta) = \frac{1}{|D|} \sum_{(x,y) \in D} l(x,y; \theta)$$

$$\theta \leftarrow \theta - \eta \underbrace{\nabla_\theta L(D; \theta)}_{\text{gradient}} = \theta - \frac{\eta}{|D|} \sum_{x,y \in D} \nabla_\theta l(x,y; \theta)$$

for $(x,y) \in D$, $x, y \sim P_D$

$M \subset D$    $|M| \ll |D|$    (minibatch)

$$\nabla_\theta L(D; \theta) \approx \mathbb{E}_{x,y \sim P_D}[\nabla l(x,y; \theta)] \approx \frac{1}{|M|} \sum_{(x,y) \in M} \nabla l(x,y; \theta)$$

Stochastic Gradient Descent

## Model selection

= hyperparameter selection / optimization

$\lambda \in \Lambda$ : a hyperparameter

$$\min_{\lambda \in \Lambda} \mathbb{E}_{(x,y) \sim D'}\left[ l'\left(x,y \; ; \; \arg\min_\theta \mathbb{E}_{(x,y) \sim D}[l(x,y; \theta)] \right) \right]$$

validation set $D'$

trng set $D$

## Reporting