## 4.6  Loss and cost function

### 4.6.1  Regression

$$\underline{y} = \underline{f}(\underline{x}; \underline{\theta}) + \underline{z} \qquad \underline{z} \text{ noise}$$

Case A:  $\underline{z} \sim N(\underline{0}, \sigma^2 \underline{\underline{I}})$ ; White Gaussian Noise

$$\underline{y} \sim N(\underline{f}(\underline{x}; \underline{\theta}), \sigma^2 \underline{\underline{I}})$$

$$q(\underline{y} | \underline{x}; \underline{\theta}) = \frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{1}{2\sigma^2} \|\underline{y} - \underline{f}(\underline{x}; \underline{\theta})\|^2\right)$$

$$l(\underline{x}, \underline{y}; \underline{\theta}) = -\ln q = \text{const} + \frac{1}{2\sigma^2} \|\underline{y} - \underline{f}(\underline{x}; \underline{\theta})\|^2$$

$$L(\underline{\theta}) = \frac{1}{N} \sum_{n=1}^{N} l(\underline{x}(n), y(n); \underline{\theta})$$

$$= \frac{1}{N} \sum_{r=1}^{N} \|y(n) - \underline{f}(\underline{x}(n); \underline{\theta})\|^2$$

mean squared error (MSE) loss, L2 loss,

method of least squares

SASP

- If $f()$ is linear in $\underline{\theta}$ (eg in perceptron)
  $\Rightarrow$ closed form solution for $\underline{\theta}$

- $f()$ for DNN → non linear in $\underline{\theta}$
  $\Rightarrow$ need iterative method  ch 4.7, Ch 5.

<u>Case B</u>  $\underline{z} \sim N(\underline{0}, \underline{\underline{C}})$ , coloured Gaussian noise

$$L(\underline{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \left( \underline{y}(n) - f(\underline{x}(n); \underline{\theta}) \right)^T C^{-1} \left( \underline{y}(n) - f(\underline{x}(n); \underline{\theta}) \right)$$

Weighted mean squared error loss.

Rarely used in practice
- How to know $\underline{\underline{C}}$ ?
- $\underline{\underline{C}}^{-1}$ is computationally very expensive.

<u>Slide 4-14</u>

4-14.1

<u>4.6.2  Classification</u>

$\underline{x} \in \mathbb{R}^d$ : input

$\underline{y} \in \{\underline{e}_1, \dots \underline{e}_c\}$ : class label for $\underline{x}$ in one-hot coding

$P_i = P(\underline{y} = \underline{e}_i \mid \underline{x})$ : true posterior

<u>ch 3.2</u> : $P(\underline{y} \mid \underline{x})$

$\qquad = \prod_{i=1}^{c} P_i^{y_i}$ : true PMF

But $P_i$ is unknown

DNN
- output $f(\underline{x}; \underline{\theta}) = [f_i(\underline{x}; \underline{\theta})] \in \mathbb{R}^c$ as estimates for $[P_i]$
- i.e. $P(\underline{y} \mid \underline{x})$ approximated by $Q(\underline{y} \mid \underline{x}; \underline{\theta})$

$\qquad\qquad\qquad = \prod_{i=1}^{c} f_i(x_i; \underline{\theta})^{y_i}$

In order to ensure

- $0 < f_i(\underline{x}; \underline{\theta}) < 1 \quad \forall i$

- $\sum_{i=1}^{c} f_i(\underline{x}; \underline{\theta}) = 1$

softmax is used in the output layer :

$$\underline{x}_L = f(\underline{x}; \underline{\theta}) = \text{softmax}(\underline{a}_L) \qquad \text{see ch 4.4}$$

$$\Rightarrow \text{loss } \ell(\underline{x}, \underline{y}; \underline{\theta}) = -\ln Q$$

$$= -\sum_{i=1}^{c} y_i \ln f_i(\underline{x}; \underline{\theta})$$

$$= -\underline{y}^T \ln f(\underline{x}; \underline{\theta}) > 0$$

Cost function $L(\underline{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \left[ -\underline{y}^T(n) \cdot \ln f(\underline{x}(n); \theta) \right]$

<u>Categorical cross entropy loss</u>

<u>Special case</u> : Binary classification, $c = 2$

ch 4.4 : softmax for $c = 2$

$\Downarrow$

one o/p neuron with sigmoid activation $f$

$$f(\underline{x}; \underline{\theta}) = \sigma(a_L)$$

Let $y_1 = y \qquad y_2 = 1 - y \quad ; \quad f_1 = f, f_2 = 1 - f$

$$\Rightarrow \ell(\underline{x}, \underline{y}; \underline{\theta}) = - \left\{ y \ln f(\underline{x}; \underline{\theta}) + (1-y) \ln(f(\underline{x}; \underline{\theta})) \right]$$

<u>binary Cross Entropy loss</u>

# 4.6.3 Semantic Image Segmentation

$$\hat{=} \text{ pixelwise classification}$$

Categorical cross entropy loss

$$\ell(\underline{\underline{X}}, \underline{Y}; \underline{\Theta}) = \sum_{h=1}^{H} \sum_{w=1}^{W} \underbrace{-\underline{y}_{hw}^{T} \ln \underline{\hat{y}}_{hw}(\underline{\underline{X}}; \underline{\Theta})}_{\text{loss for 1 pixel}}$$

$$\underbrace{\phantom{\sum_{h=1}^{H} \sum_{w=1}^{W} -\underline{y}_{hw}^{T} \ln \underline{\hat{y}}_{hw}(\underline{\underline{X}}; \underline{\Theta})}}_{\text{loss for 1 image}}$$

Problem : imbalanced classes

eg. 90% background pixels, 10% object pixels

⇒ loss cares more about the majority class

⇒ loss cares less about the minor class

⇒ reduced segmentation accuracy for minor class

Solutions :

1. Weighted categorical CE loss
2. Region-based loss

→ Let $A$ be true pos$^n$ of color pixels and $B$ be the predicted true pos$^n$ of color pixels

J and D : are not suitable for training

- ratio of integers – not differentiable
- $\underline{A}, \underline{B}, \underline{y}_{hw}$ : hard labels $\in \{0, 1\}$

  $\hat{\underline{y}}_{hw}$ : soft output / probabilities $\in R$  $0 \leq \leq 1$

Hence use Soft J and D loss

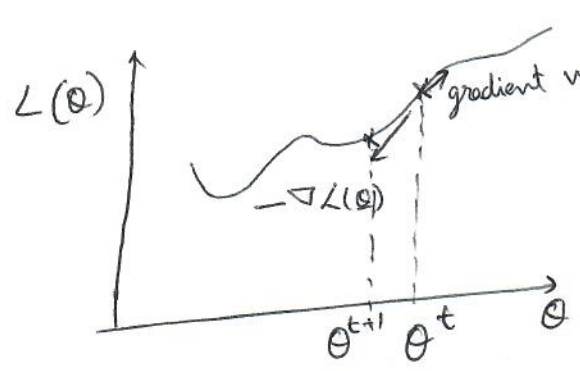<u>Slide 4-21</u>

Q for Ilias → How is soft J and D differentiable ? They still are ratio b/w 2 numbers.

## 4.7 Training

- ∘) training set $D_{train} = \{\underline{x}(n), \underline{y}(n), 1 \leq n \leq N\}$
- ∘) cost function $L(\underline{\Theta}) = \frac{1}{N} \sum_{n=1}^{N} l(\underline{x}(n), \underline{y}(n); \underline{\Theta})$
- ∘) task : $\min_{\underline{\Theta}} L(\underline{\Theta})$
- ∘) optimizer: optimization algorithm to min $L(\underline{\Theta})$

  DL : gradient descent (and variants)

  need only 1st order derivative of $L(\underline{\Theta})$



Update rule:
$$\underline{\Theta}^{t+1} = \underline{\Theta}^t - \gamma^t \nabla L(\underline{\Theta} | \underline{\Theta} = \underline{\Theta}^t)$$
↑ step size

$t = 0, \ldots$ iteration index

$\gamma^t > 0$: step size, learning rate

Calculation of $\underline{J}L(\underline{\Theta})$: non-trivial

### Chain rule of derivative:

$$\frac{d}{d\Theta} f(g(\Theta)) = \frac{df}{dg} \cdot \frac{dg}{d\Theta}$$

Layer $L$:

$$\frac{dL(\underline{\Theta})}{dW_{L,ij}} = \frac{dL(\underline{\Theta})}{d x_L} \cdot \frac{d x_L}{d a_L} \cdot \frac{d a_L}{dW_{L,ij}}$$

$$= \underline{J}_L(x_L) \cdot \underline{J}_{x_L}(\underline{a}_L) \, \underline{J}_{\underline{a}_L}(\underline{\omega}_{L,ij})$$

$$1 \times 1 \qquad = \quad \underbrace{1 \times M_L \qquad M_L \times M_L}_{\underline{J}_L(\underline{a}_L) = 1 \times M_L} \qquad M_L \times 1$$

Notation $= \underline{J}_y(\underline{x}) = \frac{dy}{d\underline{x}}$