

Q1 [10 marks]. (3 marks) Give an example of machine learning models, where the training loss is different from the error measure (i.e., measure of success). **(3 marks)** Why is such a difference desired? **(4 marks)** Should the validation criterion be the loss or error measure? And why?

Training loss: cross-entropy
 error measure: accuracy
 Advantage: differentiable training.
 Val criterion: error measure

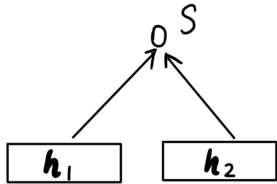
Other solutions are also acceptable:
 Training loss: MSE + ℓ_2 -penalty
 Error: MSE
 Advantage: regularization

Q2 [10 marks]. a) Show an example of two models, where the one with more parameters is more overfitting. b) Show another example, where the one with more parameters is **not** more overfitting.

a) NN with 10 hidden units VS NN with 100 hidden units
 (more parameters, more overfitting)

b) logistic regression VS 2-way softmax
 (more parameters, equivalent model)

Q3 [10 marks]. Consider a local structure of some neural networks: $s = \mathbf{h}_1^\top \mathbf{W} \mathbf{h}_2$, where $\mathbf{h}_1, \mathbf{h}_2 \in \mathbb{R}^d$, $s \in \mathbb{R}$, and \mathbf{W} is the parameter matrix.



a) **[2 marks]** Give the dimension of \mathbf{W} .

b) **[5 marks]** Assume $\partial J / \partial s$ is known. Give the recursion formulas for backpropagation through the local structure.

c) **[3 marks]** Derive the partial derivative of J with respect to \mathbf{W} .

Hints: You may use either 1) scalar calculus and organize partial derivatives in the vector form, or 2) the matrix calculus identities

- If \mathbf{A} is not a function of \mathbf{x} , then $\nabla_{\mathbf{x}} \mathbf{A} \mathbf{x} = \mathbf{A}^\top$
- If \mathbf{A} is not a function of \mathbf{x} and \mathbf{A} is symmetric, then $\nabla_{\mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$

a) $d \times d$

b) Notice that $s \in \mathbb{R}$, $\frac{\partial J}{\partial \mathbf{h}_1} = \frac{\partial J}{\partial s} \cdot \frac{\partial s}{\partial \mathbf{h}_1} = \frac{\partial J}{\partial s} \cdot \mathbf{W} \mathbf{h}_2$

$$\frac{\partial J}{\partial \mathbf{h}_2} = \frac{\partial J}{\partial s} \cdot \frac{\partial s}{\partial \mathbf{h}_2} = \frac{\partial J}{\partial s} \cdot \mathbf{W}^\top \mathbf{h}_1$$

c) $\frac{\partial s}{\partial w_{ij}} = h_{1,i} h_{2,j}$
 Thus $\frac{\partial s}{\partial \mathbf{W}} = \mathbf{h}_1 \mathbf{h}_2^\top$

$$\frac{\partial J}{\partial \mathbf{W}} = \frac{\partial J}{\partial s} \cdot \frac{\partial s}{\partial \mathbf{W}} = \frac{\partial J}{\partial s} \mathbf{h}_1 \mathbf{h}_2^\top$$

$\left[\mathbf{h}_1^\top \right] \left[-\mathbf{W} \frac{\partial J}{\partial s} \right] \left[\mathbf{h}_2 \right]$
 Printed on both sides

Q4 [25 marks]. Consider a K -way softmax classification based on input $\mathbf{x} \in \mathbb{R}^d$.

- a) [10 marks] Write the formula of the softmax classifier, and give the cross-entropy loss function.
 b) [15 marks] Show that the optimization is convex.

Hint: The bias term may be omitted for simplicity, as it may be absorbed into weights by introducing a constant feature.

$$a) \quad y_i = \text{softmax}(W\mathbf{x}) = \frac{\exp(w_i^T \mathbf{x})}{\sum_{i'} \exp(w_{i'}^T \mathbf{x})} \quad \frac{\partial J}{\partial w_{ij}} = -t_i x_j + \frac{1}{\sum_{k'} \exp(w_{k'}^T \mathbf{x})} \cdot \exp(w_i^T \mathbf{x}) \cdot x_j$$

$$\text{Loss: } J = - \sum_{k=1}^K t_k \log y_k = (y_i - t_i) x_j$$

for t_1, \dots, t_K being one-hot target

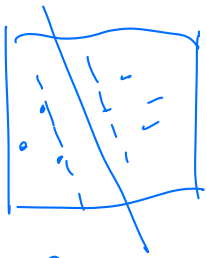
$$b) \quad J = - \sum_{k=1}^K t_k \log y_k = - \sum_{k=1}^K t_k \log \frac{\exp(w_k^T \mathbf{x})}{\sum_{k'} \exp(w_{k'}^T \mathbf{x})}$$

$$= \sum_{k=1}^K t_k \left[-w_k^T \mathbf{x} + \log \sum_{k'} \exp(w_{k'}^T \mathbf{x}) \right] = \left(- \sum_{k=1}^K t_k w_k^T \mathbf{x} \right) + \left(\log \sum_{k'} \exp(w_{k'}^T \mathbf{x}) \right)$$

Q5 [20 marks]. Consider a linearly separable binary classification, where the input is $\mathbf{x} \in \mathbb{R}^d$ and the target is $t \in \{+1, -1\}$.

- a) [10 marks] **Explain** in text the intuition of a support vector machine (SVM) and **formulate** the intuition in math. *Hint:* the distance between a point \mathbf{x}_* and a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ is $|\mathbf{w}^T \mathbf{x}_* + b| / \|\mathbf{w}\|$.
 b) [10 marks] Transform the problem into the canonical form of convex optimization. Give brief derivation steps and explanations. *Hint:* Cheatsheet in page 3. Solving SVM is not needed.

a) Maximize the margin (minimal distance between a sample and the decision boundary)



$$\text{maximize}_{w, b} \min_{m=1}^M \frac{t^{(m)} (w^T x^{(m)} + b)}{\|w\|}$$

Since w and b may scale freely, we may set $\min_{m=1}^M t^{(m)} (w^T x^{(m)} + b) = 1$

b) Then the problem becomes

$$\text{maximize}_{w, b} \frac{1}{\|w\|}$$

$$\text{s.t. } \min_{m=1}^M t^{(m)} (w^T x^{(m)} + b) = 1$$

which is equivalent to

$$\left[\begin{array}{l} \text{minimize}_{w, b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } 1 - t^{(m)} (w^T x^{(m)} + b) \leq 0 \\ \text{for } m = 1, \dots, M. \end{array} \right.$$

Q6 [25 marks]. Consider the parameters $\pi = (\pi_1, \dots, \pi_K)$ for a K -way categorical distribution.

- a) [5 marks] Given a dataset $\mathcal{D} = \{x^{(m)}\}_{m=1}^M$, where $x^{(m)} \in \{1, \dots, K\}$ is iid sampled from the categorical distribution $\text{Cat}(\pi_1, \dots, \pi_K)$. Write out the likelihood of $\pi = (\pi_1, \dots, \pi_K)$.

$$\mathcal{L}(\pi; \mathcal{D}) = p(\mathcal{D}; \pi) = \prod_{m=1}^M \prod_{k=1}^K \pi_k^{\mathbb{1}\{x^{(m)}=k\}} = \prod_{m=1}^M \prod_{k=1}^K \pi_k^{x_k^{(m)}}$$

where we define $x_k^{(m)} = \mathbb{1}\{x^{(m)}=k\}$ for simplicity.

- b) [10 marks] A Dirichlet distribution is parameterized by $\alpha = (\alpha_1, \dots, \alpha_K)$, where $\alpha_k > 0$ for $k = 1, \dots, K$. The density of a Dirichlet distribution is $p(\pi; \alpha) = C(\alpha) \prod_{k=1}^K \pi_k^{\alpha_k - 1}$, where $C(\alpha)$ is a normalizing constant depending on α .

Now suppose the prior distribution of π is a Dirichlet distribution parameterized by α . Given $\mathcal{D} = \{x^{(m)}\}_{m=1}^M$, what is the posterior distribution of π ? *Hint:* The posterior distribution would be another Dirichlet distribution parameterized by $\tilde{\alpha}$. What is $\tilde{\alpha}$?

$$\begin{aligned} \text{Posterior } p(\pi | \mathcal{D}) &\propto p(\pi) \cdot p(\mathcal{D} | \pi) \propto \prod_{k=1}^K \pi_k^{\alpha_k - 1} \cdot \prod_{k=1}^K \pi_k^{\sum_{m=1}^M x_k^{(m)}} \\ &= \prod_{k=1}^K \pi_k^{\alpha_k + \sum_{m=1}^M x_k^{(m)} - 1} \end{aligned}$$

$$\text{Thus } \pi | \mathcal{D} \sim \text{Dir}(\alpha_1 + M_1, \dots, \alpha_K + M_K) \text{ where } M_k = \sum_{m=1}^M x_k^{(m)}$$

- c) [10 marks] Derive the max a posteriori (MAP) estimation of π . *Hint:* Notice the constraint $\pi_1 + \dots + \pi_K = 1$. The constraints of $\pi_k > 0$ (for $k = 1, \dots, K$) will be automatically satisfied and thus can be ignored. (hg)

$$\text{Define } N_k = \alpha_k + \sum_{m=1}^M x_k^{(m)} - 1$$

$$\text{Goal: maximize}_{\pi} \prod_{k=1}^K \pi_k^{N_k}$$

$$\Leftrightarrow \text{minimize}_{\pi} - \sum_{k=1}^K N_k \log \pi_k$$

$$\text{Constraints: } \pi_1 + \dots + \pi_K = 1$$

Lagrangian:

$$\mathcal{L} = - \sum_{k=1}^K N_k \log \pi_k + \lambda (\pi_1 + \dots + \pi_K - 1) = 0$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = - \frac{N_k}{\pi_k} + \lambda \stackrel{\text{set}}{=} 0$$

$$\text{Thus } \pi_k = - \frac{N_k}{\lambda}$$

$$\text{Since } \pi_1 + \dots + \pi_K = 1$$

$$- \frac{N_1}{\lambda} - \dots - \frac{N_K}{\lambda} = 1$$

$$\lambda = - (N_1 + \dots + N_K)$$

$$\begin{aligned} \pi_k &= - \frac{N_k}{\lambda} \\ &= \frac{N_k}{N_1 + \dots + N_K} \\ &= \frac{\alpha_k + \sum_{m=1}^M x_k^{(m)} - 1}{\sum_{k=1}^K \alpha_k + M - K} \end{aligned}$$

END

Cheatsheet: For a convex optimization

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0 \quad \text{for } i = 1, \dots, m \\ & && h_i(\mathbf{x}) = 0 \quad \text{for } i = 1, \dots, n \end{aligned}$$

where f_i is a convex function and h_i is an affine function. The Lagrangian is defined to be

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^n \nu_i h_i(\mathbf{x})$$

For a differentiable convex optimization problem, the sufficient and necessary conditions for the optimality are

- 1) $f_i(\mathbf{x}) \leq 0$ for $i = 1, \dots, m$
- 2) $h_i(\mathbf{x}) = 0$ for $i = 1, \dots, n$
- 3) $\lambda_i \geq 0$ for $i = 1, \dots, m$
- 4) $\lambda_i f_i(\mathbf{x}) = 0$ for $i = 1, \dots, m$
- 5) $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = 0$

Scrap paper. May be detached. Additional paper is available upon request as appropriate.

May be used as answer sheets if you

- Print your name and ID on every answer sheet (including additional sheets) submitted (1 bonus mark)
- Mark your solution and corresponding problem number clearly
- Submit the sheet by the end of the exam